# A Comparative Study of Deep Reinforcement Learning Models: DQN vs PPO vs A2C

Neil de la Fuente[1,2] and Daniel Vidal[1]

[1]Autonomous University of Barcelona
[2]Computer Vision Center

December 17, 2023

## Abstract

This study conducts a comparative analysis of three advanced Deep Reinforcement Learning models – Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Advantage Actor-Critic (A2C) – exclusively within the BreakOut Atari game environment. Our research aims to assess the performance and effectiveness of these models in a singular, controlled setting. Through rigorous experimentation, we examine each model's learning efficiency, strategy development, and adaptability under the game's dynamic conditions. The findings provide critical insights into the practical applications of these models in game-based learning environments and contribute to the broader understanding of their capabilities in specific, focused scenarios.

## 1 Introduction

In this comprehensive study, we delve into the intricate world of Deep Reinforcement Learning (DRL), focusing on a thorough comparison of three renowned models: Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Advantage Actor-Critic (A2C), specifically within the BreakOut Atari game environment. To ensure consistency and robustness in our experiments, we exclusively utilized the well-established implementations of DQN, PPO, and A2C from Stable Baselines3 (SB3). This decision was driven by the desire to eliminate variability in implementation and to rely on the proven efficacy of these models in DRL applications. Our approach, centered on the standardized frameworks provided by SB3, enables a focused examination of the models' performance characteristics, free from the confounding factors that can arise from disparate model implementations. This methodology provides a clear and equitable platform for comparing the nuances of each model's learning strategy, adaptability, and efficiency within a controlled environment.

Our methodology encompassed an extensive exploration of hyperparameter settings for each model to comprehend their impact on performance in the BreakOut Atari game environment. Specifically, we systematically varied the learning rates and gamma discount factors across the DQN, PPO, and A2C models from Stable Baselines3. We experimented with some distinct learning rates for each model, aiming to identify how these rates influenced the speed and efficiency of learning, as well as the overall strategy development within the game and, we also examined the influence of two different gamma discount factors to understand how short-term versus long-term reward prioritization affected the models' decision-making processes and overall performance. The selected learning rates ranged from conservative to aggressive, providing a broad spectrum to assess each model's response to different degrees of learning aggressiveness.

This aspect of our study was pivotal in revealing each model's capability to balance immediate gratification with future gains, a key consideration in many real-world applications of DRL. Through this comprehensive approach, our research provides insightful perspectives on the optimal configuration of these models for efficient and effective learning in complex environments.

# 2 Related Work

**Deep Reinforcement Learning (DRL)** has emerged as a significant field combining deep learning with reinforcement learning, notably impacting areas like gaming and robotics. The development of DQNs by *Mnih* et al. in their pioneering work *"Playing Atari with Deep Reinforcement Learning" (Nature, 2015)* marked a turning point in DRL. This study demonstrated the effective use of neural networks with Q-learning to play Atari 2600 video games, setting a new precedent for applying deep learning to approximate Q-values. Subsequent enhancements, such as Double DQN *(van Hasselt* et al., 2016) and Dueling DQN (*Wang* et al., 2016), addressed key issues like Q-value overestimation and improved DQN's stability.

**Proximal Policy Optimization (PPO)**, introduced by *Schulman* et al. in *"Proximal Policy Optimization Algorithms" (2017)*, advanced policy gradient methods. PPO is especially noted for its efficient balance between sample efficiency and implementation ease. Its stable and reliable policy updates have made it a preferred choice in the DRL community.

**Advantage Actor-Critic (A2C)**, a simplified version of the **Asynchronous Advantage Actor-Critic (A3C)**, has been influential in actor-critic methods. A2C, while maintaining the dual advantages of learning policy and value functions, removes the complexity of asynchronous operations, as described in *Mnih* et al.'s *"Asynchronous Methods for Deep Reinforcement Learning" (2016)*.

Comparative studies in gaming environments, like the one conducted by *Henderson* et al. in *"Deep Reinforcement Learning that Matters" (2018)*, provide insights into the practical applications of these models and help in understanding their strengths and weaknesses. However, challenges like sample efficiency, stability, and generalization still persist in DRL. Current research efforts in this field are focused on addressing these issues to enhance the efficiency, scalability, and broader applicability of DRL models, as discussed in recent reviews and studies.

# 3 Methodology

In this study, we evaluated the performance of **Deep Q-Networks** (DQN), **Proximal Policy Optimization** (PPO), and **Advantage Actor-Critic** (A2C) from the **SB3** framework trained to solve the **Breakout** environment from the gymnasium library. Our experiments focused on the influence of varying hyperparameters, specifically learning rates and gamma discount factors. We tested four to five learning rates for each model to observe the effects on convergence and training robustness, and two gamma values, 0.99 and 0.90, to explore the prioritization of long-term versus immediate rewards. The training for each model was consistently conducted over a fixed number of episodes, with episode reward and learning stability as our primary performance metrics using Cuda GPU device and torch tensors to parallelize the computations and generate faster trainings. This approach facilitated a controlled and replicable comparison across the models, ensuring that any observed differences in performance were attributable to the models' intrinsic characteristics and the selected hyperparameters.

## 3.1 Implementation of the models using Stable Baselines3

Opting for SB3 allowed us to ensure that each of the three models adhered to the highest standards and best practices within the field of Deep Reinforcement Learning.

**The DQN implementation in SB3** is directly inspired by the architecture outlined in *Mnih* et al.'s landmark paper *"Playing Atari with Deep Reinforcement Learning" (Nature, 2015)*. This choice provided us with a solid, tried-and-tested baseline for DQN, replicating the neural network structure used in their pioneering study. The consistent and reliable performance of the SB3 version of DQN offered a robust foundation for our comparative exploration, particularly within the context of the BreakOut environment.

For the **PPO and A2C models**, we also relied on the robust implementations available in SB3. PPO is particularly notable for its balance between sample efficiency and simplicity in implementation, aligning with the original design principles laid out by *Schulman* et al. in *"Proximal Policy Optimization Algorithms" (2017)*. The A2C model in SB3 maintains the core advantages of A3C while removing the need for asynchronous operations, as detailed in *Mnih* et al.'s *"Asynchronous Methods for Deep Reinforcement Learning" (2016)*.

By leveraging these established and validated models from SB3, our study benefits from a high level of consistency and comparability. This approach allows us to focus intently on analyzing how each model

performs in the BreakOut environment under various hyperparameter configurations.

## 3.2 Hyperparameter Variations

In order to evaluate the impact on the performance of the different models, we explored a spectrum of learning rates to identify how different rates affected each model's speed of learning and robustness of strategy development. In tandem, we manipulated the gamma discount factor with two distinct values, 0.99 and 0.90, to assess the models' sensitivity to short-term versus long-term rewards. This experimental design was aimed at discovering configurations that optimized learning efficiency while being mindful of computational resource constraints. Our structured approach provides a foundation for an in-depth evaluation of the learning dynamics and performance efficacy of each model, which we will further elucidate in the 'Experiment Setup' and 'Results' sections.

# 4 Hypotheses and Theoretical Considerations

In anticipation of our experimental results, we formed several hypotheses based on the theoretical underpinnings of the models in question and their known strengths and weaknesses. We titled this section "Hypotheses and Theoretical Considerations" to reflect both the predictive and the conceptual nature of the content. This section is structured to delineate our theoretical expectations into distinct subsections, each corresponding to a core hypothesis about the performance of the models in the specific setting of the BreakOut Atari game.

## 4.1 Alignment of BreakOut Dynamics with DQN's Value Estimation

In the BreakOut Atari game, the dynamics are particularly well suited for the Q-learning algorithm, upon which the Deep Q-Network (DQN) is based. Q-learning seeks to learn a value function $Q(s, a)$, representing the expected return of taking an action $a$ in a state $s$ and following the optimal policy thereafter. For BreakOut, the state $s$ can be described by the position of the paddle, the ball, and the configuration of the bricks, while the action $a$ corresponds to moving the paddle left, right, or staying in place.

The game's immediate and clear reward structure—points gained from breaking bricks—aligns with the value estimation approach of DQN. The Q-function for DQN is updated using the following rule:

$$Q_{new}(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]$$

where $\alpha$ is the learning rate, $\gamma$ is the discount factor, $r_{t+1}$ is the immediate reward, and $\max_{a'} Q(s_{t+1}, a')$ is the maximum predicted value for the next state $s_{t+1}$, over all possible actions $a'$. This formula allows DQN to incrementally improve its policy by learning from the immediate outcomes of its actions, which in BreakOut are straightforward and tightly coupled with the action taken.

Contrastingly, on-policy methods like PPO and A2C involve estimating the policy directly. The policy $\pi(a|s)$ represents the probability of taking action $a$ in state $s$, and is adjusted in the direction suggested by the policy gradient:

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

The advantage function $A^{\pi_\theta}(s_t, a_t)$, which measures the relative value of an action compared to the average, is more complex to estimate in a game like BreakOut. The reason is that certain strategies, like "tunneling" where the ball is directed behind the brick wall to clear multiple bricks, require foresight and a nuanced understanding of the game's physics, which are not immediately obvious from the reward signal.

Estimating the value of such strategies is more challenging than directly mapping actions to outcomes, as required in Q-learning. DQN's value estimation through Q-learning is more straightforward because it can directly correlate actions like "move paddle under the ball" with receiving points, whereas PPO and A2C must learn through trial and error which complex sequences of actions yield higher returns. This makes DQN more suited for BreakOut, where the optimal policy is easier to learn through the value function rather than direct policy estimation.

## 4.2 Advantage of Experience Replay in DQN

The concept of experience replay is fundamental to the success of DQN. In the context of the BreakOut game, the mechanics often lead to repetitive scenarios where the paddle must hit the ball in a rhythmical pattern. These recurrent interactions between the ball and paddle could potentially lead to a form

of learning that is too specialized—optimized for the frequent but limited scenarios encountered.

The experience replay mechanism mitigates this risk by storing a history of experiences:

$$e_t = (s_t, a_t, r_{t+1}, s_{t+1})$$

in a replay buffer $D_t$. This buffer acts as a diversified data reservoir from which the DQN samples to update the value function.

This process ensures that the DQN does not simply memorize the most recent or frequent patterns but instead develops a more generalized understanding of the game by learning from a broader range of experiences.

The advantage of such an approach in BreakOut is twofold. First, it allows the DQN to break free from the constraints of local optima—situations where a certain repetitive strategy might seem best because it has been the most experienced so far, but is not globally optimal. Second, it enables the DQN to learn from rare but informative episodes that could be crucial for mastering the game. For example, if the ball reaches the back of the brick wall creating a 'tunnel', which is a less frequent but highly rewarding event, the replay buffer ensures that these valuable experiences are not lost.

By continuously revisiting a diverse set of past experiences, the DQN's training process becomes more robust and less prone to overfitting. This is particularly beneficial in games like BreakOut, where the difference between a good and a great strategy can often be subtle, and learning from the full breadth of past experiences is key to discovering these nuances. Experience replay, therefore, not only enhances the learning of effective strategies but also the retention of them, which is essential for mastering environments with repetitive and pattern-based dynamics.

## 4.3 Sensitivity to Hyperparameters Across Models

The sensitivity of an algorithm to its hyperparameters can greatly influence its performance, especially in complex and dynamic environments such as the BreakOut Atari game. Our hypothesis posits that DQN and PPO will be less sensitive to learning rate variations, while A2C will display greater sensitivity.

DQN's use of experience replay and fixed Q-targets provides a stabilizing effect on learning updates. The experience replay randomizes the data, thus breaking correlations in the observation sequence and smoothing changes in the data distribution. This leads to more stable gradient descent updates, allowing DQN to handle a wider range of learning rates.

The robustness of PPO comes from its objective function, which uses a clipping mechanism to prevent excessively large policy updates that could lead to instability. This clipping mechanism acts as a safeguard against the potential negative effects of choosing a suboptimal learning rate.

On the other hand, A2C's performance is more intricately tied to its learning rate. As an on-policy algorithm, A2C continuously updates its policy based on the latest data it collects from the environment. If the learning rate is too low, A2C's policy may not adapt quickly enough to the new information, leading to suboptimal performance. Conversely, a high learning rate might cause the policy to change too rapidly, which can destabilize the learning process. The ideal learning rate for A2C should therefore strike a balance, allowing it to quickly integrate new data without causing instability.

Discount factors also play a crucial role in the learning process. A high discount factor ($\gamma$) makes an algorithm more farsighted by placing more emphasis on future rewards. This could potentially be more beneficial for DQN, as it may enable the agent to better recognize the long-term benefits of strategic moves like tunneling. A lower $\gamma$ could make the agent myopic, prioritizing immediate rewards, which might suffice for simpler strategies but fail to capture the depth of the game's strategic possibilities.

In the case of PPO and A2C, a well-chosen discount factor helps in balancing the trade-off between exploring new strategies and exploiting known rewarding behaviors. PPO's ability to maintain a stable update might be less affected by the discount factor compared to A2C, as the latter can be more sensitive to the choice of $\gamma$ due to its direct policy update mechanism.

In summary, we hypothesize that DQN and PPO will exhibit a higher tolerance to hyperparameter variations, while A2C will require careful tuning of its learning rate and discount factor to achieve optimal performance. This sensitivity to hyperparameters is an essential consideration in the deployment and tuning of reinforcement learning models in practice.

## 4.4 Policy Optimization Stability in PPO and A2C

The stability of policy optimization in Proximal Policy Optimization (PPO) and Advantage Actor-Critic (A2C) hinges on their distinct approaches to updating policies. PPO is designed to mitigate policy volatility through a clipping mechanism in its objective function, which constrains the extent of policy updates, fostering stability. The clipped objective in PPO is given by:

$$\mathrm{L}^{CLIP}(\theta) = E_t \left[ \min(r_t(\theta)\hat{A}_t, clip(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \right]$$

where $r_t(\theta)$ is the ratio of the new policy to the old policy probabilities, $\hat{A}_t$ is the estimator of the advantage function at time $t$, and $\epsilon$ is a hyperparameter that defines the clipping range.

Conversely, A2C updates policies after every step, lacking PPO's clipping guardrails, which can lead to more aggressive policy shifts:

$$\Delta\theta = \alpha\nabla_\theta \log \pi_\theta(a_t|s_t)A(s_t, a_t)$$

The absence of a clipping mechanism in A2C allows for larger, and potentially more disruptive, policy updates, especially when the advantage $A(s_t, a_t)$ is significant.

We anticipate that PPO will demonstrate more stable learning progression due to its conservative update strategy, while A2C may show more variability in its performance, potentially achieving higher peaks and experiencing deeper troughs in response to the BreakOut game's dynamic challenges. This variability is reflective of A2C's responsiveness to immediate changes in the environment, which can lead to rapid but sometimes unstable learning trajectories.

## 4.5 Conclusions on hypotheses

These hypotheses serve not only as predictions to be tested against the empirical data but also as a framework to interpret the complexities of model behaviors. They reflect a balance between established theoretical knowledge and the particularities of the BreakOut Atari environment. The subsequent sections will detail the experimental setups and results, providing a platform to assess the validity of these theoretical considerations.

# 5 Experiment Setup

To thoroughly examine the capabilities of Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Advantage Actor-Critic (A2C), our experiments were conducted within the BreakOut Atari game environment, leveraging the Gymnasium framework by the Farama Foundation.

## 5.1 Training Environment

The BreakOut Atari game provided by Gymnasium stands out as an exemplary testbed due to its historical significance in reinforcement learning and an ideal balance of complexity. It challenges the learning algorithms to master both reactive and strategic gameplay while being simple enough to ensure experimental reproducibility.

## 5.2 Parameter Variations

The experimental design hinged on the careful modulation of key hyperparameters, crucial for examining the adaptability and efficiency of the DRL models. The parameters adjusted included:

- **Learning Rate**: We experimented with learning rates at several magnitudes: $1e-5, 5e-5, 1e-4, 5e-4, 1e-3, 5e-3$. This broad spectrum allowed us to investigate the impact of learning rate adjustments on the convergence and performance stability of the models.

- **Gamma Discount Factor**: To understand the influence of temporal discounting on the strategic planning of the models, we varied the gamma discount factor between $\gamma = 0.99$ and $\gamma = 0.90$. This exploration aimed to assess how the models prioritize short-term gains versus long-term rewards.

Selecting these hyperparameters was pivotal since the learning rate controls the magnitude of updates to the model's knowledge base, and the gamma discount factor balances the emphasis on immediate versus delayed rewards. These elements are fundamental to the development of robust and effective learning strategies in reinforcement learning agents.

For our experiments, training duration was quantified in terms of frames rather than time or episodes, providing a consistent measure across all models. Each model underwent training for approximately 20 million frames. This metric was chosen as it offers an objective standard of comparison, ensuring no model is inadvertently favored by the measurement approach.

The frame-based approach aligns with the conventions of the Stable Baselines3 (SB3) framework and

allows for a fair assessment of each model's learning process, as it accounts for variations in episode lengths. For instance, more skilled agents may play longer, resulting in fewer but lengthier episodes. This is particularly relevant in environments like BreakOut, where proficient gameplay can significantly extend the duration of a single episode.

An additional consideration is the potential for the game to enter into loops, such as the ball bouncing between the wall and ceiling repeatedly without hitting bricks. Such scenarios can artificially inflate the number of frames without corresponding to meaningful learning or gameplay progress. This effect is evident in the learning curves, where certain training runs show extended plateaus or spikes in episode duration, which may indicate the agent has entered a loop rather than achieving a breakthrough in strategy.

By standardizing the training duration across models via frame count, we aimed to mitigate these factors, providing a clear, unbiased view of each model's learning efficiency and performance over the course of consistent environmental interaction.

## 5.3  Evaluation Metrics

To assess the effectiveness of each model, we employed multiple evaluation metrics that together offer a holistic view of performance. These metrics are critical as they capture various aspects of the learning process and the models' proficiency in navigating the BreakOut environment.

- **Average Reward per Episode**: This metric measures the mean score obtained by the agent per episode over the course of training. A higher average reward indicates a more successful strategy in maximizing points, which in BreakOut, correlates with the agent's ability to keep the ball in play and efficiently break bricks, noting that the reward of each brick is dependant on its color.

- **Episodes to Threshold**: We tracked the number of episodes each model needed to consistently reach predefined reward thresholds. This metric sheds light on the learning speed of each model—how quickly it can grasp and apply effective strategies within the game.

- **Time to Threshold**: We measured the amount of hours each model needed to reach predefined reward thresholds. This metric sheds light on

the learning speed of each model—how quickly it can grasp and apply effective strategies within the game.

- **Reward Distribution**: Understanding the spread and variation of rewards received can indicate the consistency of the model's performance. A narrow distribution suggests stability, while a wide distribution may reflect a more explorative or unstable learning process.

- **Stability of Learning**: The fluctuation in the models' performance over time was monitored to assess their learning stability. Frequent and large variations in performance might indicate a model is struggling to converge to an optimal policy, while smooth progression suggests steady learning.

- **Frame Utilization Efficiency**: Considering the potential for the game to enter nonterminating loops, we also evaluated how efficiently models used the allocated frames. This metric highlights the models' capacity to utilize experiences for learning purposefully, avoiding wasteful or redundant gameplay.

Each of these metrics contributes to a composite picture of model performance. By examining these various facets, we can discern not only which model achieved the highest scores but also understand the nuances of their learning processes. This multifaceted evaluation is designed to dissect the models' responses to different configurations within the standardized BreakOut environment, setting the stage for a detailed analysis of their respective proficiencies and weaknesses.

# 6 Results

Our empirical investigation into the performance of Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Advantage Actor-Critic (A2C) within the BreakOut Atari game environment yielded insightful findings, as detailed in the subsections below.
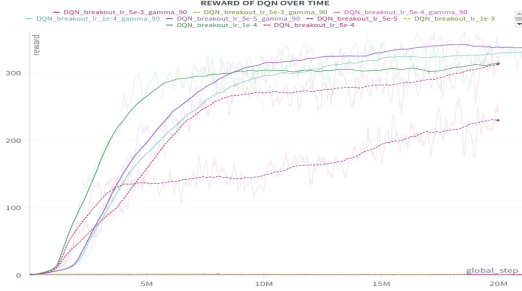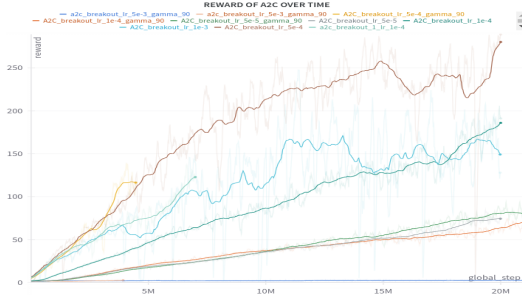
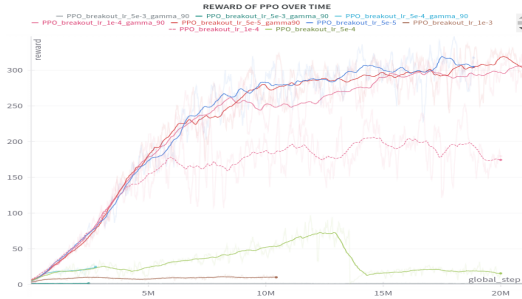

Figure 1: DQN performance



Figure 2: A2C Performance



Figure 3: PPO Performance

## 6.1 Performance Across Learning Rates

The learning rate is a pivotal hyperparameter in reinforcement learning, dictating the step size during the optimization process. From the data, DQN exhibited remarkable resilience across a broad span of learning rates, consistent with our hypothesis that its experience replay buffer would confer robustness against the variability in learning rates. Specifically, DQN main-

tained a relatively steady improvement in reward as the learning rate increased, up to a threshold beyond which performance plateaued or slightly deteriorated, indicating an optimal learning rate range.

PPO also showed resistance to changes in learning rate, though its performance was more sensitive to the extremes. At the lower end of the learning rate spectrum, PPO's progress was gradual but consistent. However, at higher rates, PPO's performance was more variable, likely due to the algorithm's clipped objective function which, while mitigating the risk of large destructive updates, also caps the potential rapid advancements that can be achieved with more aggressive learning rates.

In contrast, A2C's performance was heavily influenced by the learning rate. At lower rates, A2C's learning curve was almost dead, suggesting an inability to make significant policy improvements. As the learning rate increased, A2C's performance improved markedly, confirming our prediction that A2C requires a careful balance in learning rate settings to optimize its continuous policy updates.

## 6.2 Adaptability to Discount Factor Variations

The gamma discount factor influences the agent's consideration of future rewards. Across all models, a higher gamma value typically correlated with a more strategic play, where long-term gains were prioritized. DQN's performance peaked at a moderate gamma value before declining, suggesting a sweet spot where the agent was sufficiently farsighted without being hindered by the overvaluation of distant future rewards.

PPO and A2C demonstrated a clearer preference for a higher gamma value. This was particularly true for PPO, which displayed a steady increase in performance as the gamma value rose, aligning with the algorithm's inherent stability and its capability to evaluate long-term strategies effectively.

## 6.3 Learning Stability, Efficiency, and Reward Optimization

In examining the learning process of DQN, PPO, and A2C models, we considered several critical aspects: the stability and efficiency of learning, reward optimization, and episode length within the BreakOut Atari game. DQN's performance stood out, demonstrating a smooth learning curve and efficient frame

utilization, highlighting its capability to rapidly assimilate and apply successful strategies. This was mirrored in its real-time training efficiency, where DQN consistently achieved high rewards in shorter durations, making it an ideal candidate for scenarios demanding quick adaptation and time-efficient learning.
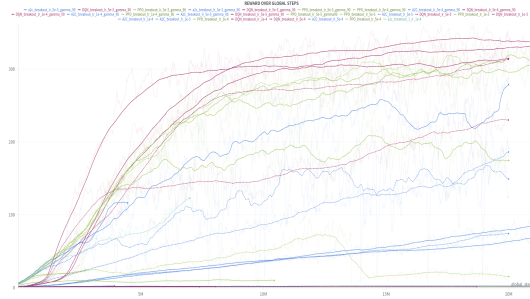


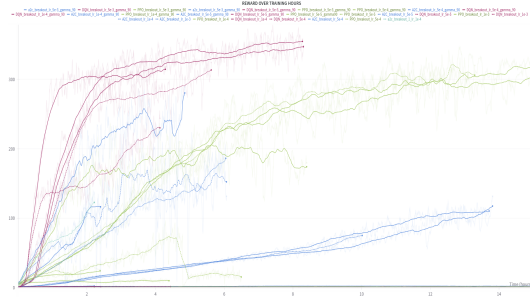Figure 4: Reward over training steps



Figure 5: Reward over training hours

Contrastingly, PPO and A2C experienced more pronounced fluctuations in their learning trajectories. PPO, though displaying a reasonable level of stability, tended to engage in lengthier episodes, suggesting a propensity for a more exploratory approach that may sacrifice swiftness for thoroughness. A2C's learning curve was the most variable, reflecting its sensitivity to environmental dynamics and possibly a greater need for exploration to refine its policy. The higher number of frames A2C required to reach performance levels comparable to DQN implies a less efficient learning process, especially marked when considering the model's real-time performance, which lagged behind the others. Furthermore, reward optimization and episode length analysis revealed strategic distinctions among the models. DQN's strategy excelled in securing high scores efficiently, translating learning into performance gains with remarkable swiftness. In contrast, PPO and A2C, despite their eventual improvements, necessitated longer episodes for similar levels of reward, indicating potentially less efficient strategies. This behavior underscores a strategic divergence where DQN prioritizes exploitation of the known dynamics, while PPO and A2C seem to invest heavily in exploration, a trait that may yield substantial benefits in less structured and more complex environments.

The practical implications of our findings are significant. They suggest that while DQN shines in environments requiring rapid learning and deployment, PPO and A2C may offer advantages in scenarios where the depth of exploration and comprehensive strategy development are key. Such insights are invaluable for the nuanced selection and application of reinforcement learning models, emphasizing the importance of aligning the model's strengths with the specific demands and constraints of the task environment.
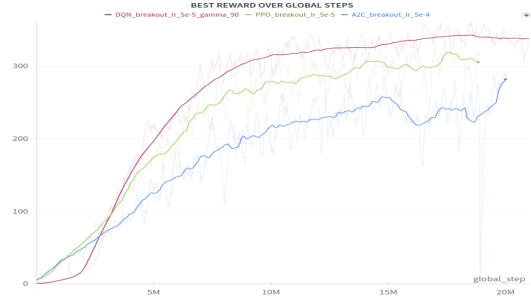


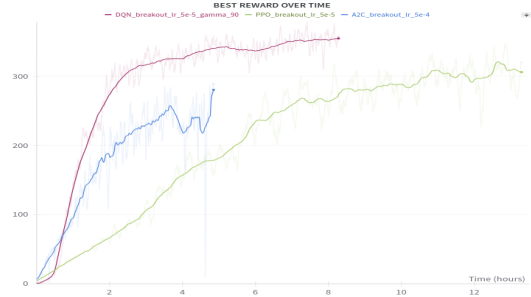Figure 6: Reward over training time for the best model of each type



Figure 7: Reward over training global time for the best Model of each type

## 6.4 Summary of results

In summary, our results substantiate the hypotheses posed in our theoretical considerations, illustrating the nuanced interplay between model architectures, hyperparameters, and the BreakOut game dynamics. The collected data provide a comprehensive picture of how each model's unique characteristics influence its learning trajectory and overall performance in a standardized environment.

# 7 Findings

Upon rigorous comparison of Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Advantage Actor-Critic (A2C) within the structured confines of the BreakOut Atari game, our study brings to light nuanced insights that prompt a reevaluation of the preferred use cases for these models in the realm of Deep Reinforcement Learning (DRL).

## 7.1 DQN: A Contender for Strategic Efficiency

The study's revelations underscore DQN's unexpected superiority in environments with clear, immediate reward structures. DQN's methodical approach to value estimation is well-suited to such scenarios, allowing for a rapid and efficient development of effective strategies. This efficiency is bolstered by the experience replay buffer, a feature that enables DQN to draw from a diverse set of past experiences, thus preventing over-specialization and fostering a robust policy that generalizes across a multitude of game states.

## 7.2 Model Adaptability and Hyperparameter Resilience

DQN's performance exhibits commendable resilience to hyperparameter fluctuations, making it a versatile and forgiving model for practitioners. This adaptability is contrasted with PPO and A2C, which, despite their advanced policy optimization capabilities, show a pronounced sensitivity to hyperparameter settings. This necessitates a fine-grained tuning process to navigate the trade-offs between exploration and exploitation, especially in environments where the reward pathways are less direct and the strategic demands are higher.

## 7.3 Strategic Exploration: PPO and A2C's Domain

In environments that reward deep exploration and complex strategy development, PPO and A2C demonstrate their prowess. Their policy gradient methods, equipped to probe the depths of a more intricate state space, become advantageous. While these models may require extended periods to converge on highly rewarding strategies, their potential in tasks demanding a sophisticated level of strategic planning is evident.

## 7.4 Towards a Contextual Model Selection Framework

The comparative performance of DQN, PPO, and A2C accentuates the necessity for a contextual approach to model selection in DRL. Our findings suggest that while DQN is optimally poised for tasks demanding quick learning and efficient adaptation, environments characterized by their opaqueness and strategic complexity may benefit from the exploratory strengths of PPO and A2C.

## 7.5 Implications for Practical Application

The implications of our findings for the practical application of DRL are multifaceted. They serve as a guide for practitioners in choosing a model that not only fits the immediate needs of the task but also aligns with the long-term objectives of the learning process. As our understanding of DRL models deepens, it becomes clear that the decision-making framework for selecting a DRL model must be as dynamic and multifaceted as the models themselves.

In closing, the study advocates for a discerning application of DRL models, where the choice is not dictated by the perceived sophistication of the algorithm but by a strategic fit for the task. This approach ensures that the selected model is not only theoretically sound but also practically effective, able to harness the unique dynamics of the environment for optimal learning and performance.

# 8 Conclusion and Future Work

This study has meticulously evaluated Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and Advantage Actor-Critic (A2C) within the BreakOut Atari environment, unveiling distinct capabilities and performance profiles for each model. DQN, traditionally viewed as a simpler model, has demonstrated its robustness and efficiency in a controlled setting, suggesting that its utility in practical applications remains significant, especially in environments with clear, immediate reward structures. PPO and A2C, while typically favored for complex tasks due to their advanced policy optimization techniques, require careful tuning and strategy development, which can be advantageous in more intricate or less predictable environments. The findings highlight

the importance of choosing the right model based on the specific characteristics and requirements of the

task environment. This study emphasizes the need for a nuanced approach to model selection that goes beyond the complexity of the algorithm to consider strategic alignment with the task at hand.

Future research should expand on the comparative analysis of DQN, PPO, and A2C across a wider range of environments, including those with delayed rewards and higher complexity. There is also an opportunity to explore the integration of hybrid models that combine the strengths of value-based and policy-based approaches. Further studies could investigate the impact of additional hyperparameters, network architectures, and reward shaping techniques on the performance of these models. Moreover, the development of adaptive algorithms that can dynamically select or switch between models based on real-time performance metrics within an environment could be highly beneficial. This would contribute to the creation of more robust and versatile DRL systems, capable of adjusting to the nuances of various tasks and maximizing learning efficiency. Lastly, the exploration of DRL applications in real-world scenarios, such as robotics, autonomous vehicles, and financial modeling, where the environment is often unpredictable and complex, will be essential in translating these findings into tangible benefits and advancements in the field of artificial intelligence.

# References

[1] Mnih, V., Kavukcuoglu, K., Silver, D., et al. *Playing Atari with Deep Reinforcement Learning.* arXiv preprint arXiv:1312.5602, 2013.

[2] Raffin, A., Hill, A., Gleave, A., et al. *Stable Baselines3: Reliable Reinforcement Learning Implementations.* Journal of Open Source Software, 2021.

[3] Mnih, V., Kavukcuoglu, K., Silver, D., et al. *Human-level control through deep reinforcement learning.* Nature, 518(7540), 529-533, 2015.

[4] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. *Proximal Policy Optimization Algorithms.* arXiv preprint arXiv:1707.06347, 2017.

[5] Mnih, V., Badia, A. P., Mirza, M., et al. *Asynchronous Methods for Deep Reinforcement Learning.* In Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48 (ICML'16). JMLR.org, 1928–1937.

[6] Henderson, P., Islam, R., Bachman, P., et al. *Deep Reinforcement Learning that Matters.* In Proceedings of the AAAI Conference on Artificial Intelligence, 2018.

[7] Van Hasselt, H., Guez, A., and Silver, D. *Deep Reinforcement Learning with Double Q-learning.* In Proceedings of the AAAI Conference on Artificial Intelligence, 2016.

[8] Wang, Z., Schaul, T., Hessel, M., et al. *Dueling Network Architectures for Deep Reinforcement Learning.* In Proceedings of the 33rd International Conference on Machine Learning, 2016.

[9] Sutton, R. S., and Barto, A. G. *Reinforcement Learning: An Introduction (Second edition).* MIT Press, 2018.

[10] Lillicrap, T. P., Hunt, J. J., Pritzel, A., et al. *Continuous control with deep reinforcement learning.* arXiv preprint arXiv:1509.02971, 2016.

[11] Silver, D., Schrittwieser, J., Simonyan, K., et al. *Mastering the game of Go without human knowledge.* Nature, 550(7676), 354-359, 2017.

[12] Hessel, M., Modayil, J., Van Hasselt, H., et al. *Rainbow: Combining Improvements in Deep Reinforcement Learning.* In Proceedings of the AAAI Conference on Artificial Intelligence, 2018.

[13] Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. *The Arcade Learning Environment: An Evaluation Platform for General Agents.* Journal of Artificial Intelligence Research, 47:253-279, 2013.

[14] Brockman, G., Cheung, V., Pettersson, L., et al. *OpenAI Gym.* arXiv preprint arXiv:1606.01540, 2016.

[15] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. *Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor.* In Proceedings of the International Conference on Machine Learning, 2018.

[16] Schaul, T., Quan, J., Antonoglou, I., and Silver, D. *Prioritized Experience Replay.* arXiv preprint arXiv:1511.05952, 2015.

[17] Duan, Y., Chen, X., Houthooft, R., Schulman, J., and Abbeel, P. *Benchmarking Deep Reinforcement Learning for Continuous Control.* In Proceedings of the International Conference on Machine Learning, 2016.

[18] Tassa, Y., Doron, Y., Muldal, A., et al. *DeepMind Control Suite.* arXiv preprint arXiv:1801.00690, 2018.

[19] Coumans, E., and Bai, Y. *PyBullet, a Python module for physics simulation for games, robotics and machine learning.* GitHub repository, 2016.

[20] The Farama Foundation *Gymnasium: A fork of OpenAI Gym with additional environments and functionalities.*