

# Bayesian Optimization of Hyperparameters for Rainbow DQN in the CartPole-v1 Environment

**Abstract**—This paper presents a Bayesian optimization approach to hyperparameter tuning for the Rainbow DQN reinforcement learning algorithm, using the Hyperopt library and the CartPole-v1 environment as a benchmark. The study investigates the impact of search space definition on the convergence and quality of optimized hyperparameters. Furthermore, it analyzes the effectiveness of different evaluation methods in the context of hyperparameter optimization for deep reinforcement learning. Results demonstrate the efficacy of Bayesian optimization in identifying high-performing hyperparameter configurations for Rainbow DQN in this control task.

## I. INTRODUCTION

Deep Reinforcement Learning (DRL) has gained substantial attention for its ability to solve complex control tasks by combining reinforcement learning with deep neural networks. Among DRL algorithms, the Deep Q-Network (DQN) has demonstrated remarkable success, notably in environments with discrete action spaces such as Atari games. However, the performance of DQN and its variants is highly sensitive to the choice of hyperparameters. Selecting optimal configurations often involves a costly and time-consuming trial-and-error process.

Rainbow DQN [1] is a notable advancement in the DQN family, integrating multiple enhancements—such as Double DQN, Prioritized Experience Replay, Dueling Networks, Noisy Nets, and Categorical DQN—into a single, unified agent. While Rainbow DQN provides superior performance across several benchmark tasks, its numerous hyperparameters introduce complexity in tuning and reproducibility. Suboptimal settings can degrade learning efficiency and final performance.

To address this challenge, this paper explores the use of Bayesian optimization for automatic hyperparameter tuning of Rainbow DQN in the widely used CartPole-v1 environment. Bayesian optimization, known for sample efficiency in high-dimensional spaces, models the objective function using probabilistic surrogates and balances exploration and exploitation to guide the search. We leverage the `Hyperopt` library with a Tree-structured Parzen Estimator (TPE) approach to evaluate the impact of various hyperparameter search spaces.

Our key contributions are as follows:

- We apply Bayesian optimization to efficiently tune the hyperparameters of Rainbow DQN for the CartPole-v1 task.
- We investigate how the definition of the search space influences convergence behavior and solution quality.

- We compare the optimized Rainbow agent against baseline configurations and ablated versions to assess the contribution of individual components.

Through systematic experimentation and analysis, we demonstrate that Bayesian optimization can effectively identify high-performing configurations for Rainbow DQN, achieving robust and stable learning outcomes in a computationally efficient manner. The insights gained from this study are valuable for practitioners seeking to deploy DRL agents with minimal manual tuning effort.

## II. RELATED WORK

Hyperparameter optimization and architectural improvements in Deep Reinforcement Learning (DRL) have been widely studied to enhance sample efficiency and policy performance. This section reviews key developments in both domains, focusing on Bayesian optimization for hyperparameter tuning and the evolution of DQN variants that culminated in Rainbow DQN.

### A. Hyperparameter Optimization in Deep RL

Hyperparameters such as learning rate, discount factor, and exploration strategies play a critical role in the stability and convergence of DRL algorithms. Manual tuning is both time-consuming and task-specific, motivating the development of automated optimization methods. Bayesian optimization is a widely adopted technique for this purpose due to its sample-efficient, model-based search strategy [2]. The Tree-structured Parzen Estimator (TPE), introduced by Bergstra et al. [3], enables flexible optimization over both discrete and continuous domains and has been effectively applied in DRL contexts [4].

Recent work has integrated Bayesian methods into DRL pipelines, including Gaussian Processes [5], surrogate-based meta-learning [6], and Monte Carlo tree-based search for multi-fidelity optimization [7]. Open-source libraries such as Hyperopt [8], Optuna [9], and BOHB [10] have accelerated experimentation in this field. Janner et al. [11] and Song et al. [12] evaluated automated hyperparameter tuning frameworks in environments like MuJoCo and Atari.

### B. DQN and Its Extensions

The original DQN architecture introduced by Mnih et al. [13] demonstrated that deep networks can approximate value functions effectively when combined with experience replay and target networks. However, it suffers from value overestimation and slow convergence, prompting several extensions:

- **Double DQN** reduces overestimation bias by decoupling action selection and value estimation [14].
- **Dueling DQN** introduces separate streams for value and advantage functions, improving learning in environments with similar action values [15].
- **Prioritized Experience Replay (PER)** samples experiences based on TD-error magnitudes, improving sample efficiency [16].
- **Noisy Networks** facilitate efficient exploration by injecting learned noise into parameters [17].
- **Distributional RL** models the entire return distribution rather than just the expectation, leading to more stable learning [18].

The Rainbow DQN architecture proposed by Hessel et al. [1] combined all these improvements, yielding superior performance across the Atari 2600 benchmark suite. Subsequent work examined Rainbow’s generalization capacity in new environments and its integration with curiosity-based exploration [19], meta-learning [20], and multi-agent systems [21].

### C. Combined Optimization and Architecture Studies

Several recent studies have combined architectural innovations with automated hyperparameter tuning. For instance, Runge et al. [22] and Zahavy et al. [23] explored adaptive agent architectures whose parameters are tuned during training. Chen et al. [24] proposed AutoRL frameworks that co-optimize neural architecture and learning settings.

In parallel, Hutter et al. [25] and Elsken et al. [26] surveyed neural architecture search (NAS) and automated machine learning (AutoML) techniques applicable to RL settings. These systems further reduce manual intervention by searching across architecture and optimization strategy spaces jointly.

### D. Implications for Lightweight Environments

Though most studies focus on complex benchmarks, tasks like CartPole remain relevant for benchmarking and rapid prototyping. Studies by Brockman et al. [27] and Henderson et al. [28] emphasized reproducibility and hyperparameter sensitivity even in simple environments, advocating for tuning reproducibility standards.

In summary, Rainbow DQN represents a well-founded combination of independently validated techniques. Its performance can be further improved with principled hyperparameter tuning via Bayesian optimization. This paper contributes to the growing literature on bridging algorithmic innovations in DRL with efficient hyperparameter search strategies.

## III. BACKGROUND

In reinforcement learning (RL), agents aim to learn optimal policies by interacting with an environment and receiving feedback in the form of rewards. Deep Q-Networks (DQNs) have been pivotal in scaling Q-learning to high-dimensional state spaces by leveraging deep neural networks. Over time, several extensions and improvements have been proposed to overcome limitations of the original DQN framework. Rainbow DQN integrates several of these enhancements into a unified

agent, including Double DQN, Prioritized Experience Replay, Dueling Networks, Noisy Nets, and Categorical DQN [1]. This section presents the conceptual background of each component forming the foundation of Rainbow DQN.

### A. Deep Q-Networks

Traditional Q-learning struggles in environments with large or continuous state spaces due to its reliance on discrete Q-tables. DQN addresses this by approximating the Q-function using a deep neural network, which takes states as input and outputs Q-values for each possible action. The agent selects actions based on an  $\epsilon$ -greedy policy and uses experience replay to break correlations between sequential data. A target network, updated less frequently, provides stable targets during training. This combination significantly improves convergence in complex environments.

### B. Double DQN

DQN often overestimates action values, leading to sub-optimal policies. Double DQN mitigates this by decoupling the action selection and evaluation steps. While one network chooses the best action, the target network estimates its value, thus reducing overestimation bias. Formally, the target value is computed as:

$$y = r + \gamma Q_{\theta^-}(s', \arg \max_a Q_{\theta}(s', a))$$

where  $\theta$  and  $\theta^-$  denote parameters of the main and target networks, respectively [14].

### C. Prioritized Experience Replay

Experience replay in standard DQN samples transitions uniformly, which may not emphasize critical learning experiences. Prioritized Experience Replay improves learning efficiency by sampling experiences based on the temporal-difference (TD) error. Transitions with higher TD error—indicating greater surprise or learning potential—are sampled more frequently. A stochastic prioritization mechanism, typically using a probability proportional to the TD error’s magnitude, is used to balance exploration and overfitting [16].

### D. Dueling DQN

Dueling network architecture separates the estimation of the state-value function  $V(s)$  and the advantage function  $A(s, a)$  within the neural network. This helps the network learn which states are (or are not) valuable, independently of the chosen action. The final Q-value is computed as:

$$Q(s, a) = V(s) + \left( A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a'} A(s, a') \right)$$

This formulation improves learning stability and decision-making in scenarios where some actions do not significantly impact rewards [15].

### E. Noisy Nets

Exploration is critical in reinforcement learning. Instead of using simple random strategies like  $\epsilon$ -greedy, Noisy Nets introduce learnable noise into the network weights, allowing for adaptive exploration. Parameters of the noise distribution are updated during training, enabling the agent to learn optimal exploration strategies over time. The stochasticity introduced leads to more consistent and efficient discovery of valuable state-action pairs [17].

### F. Categorical DQN

Categorical DQN represents the Q-value distribution instead of estimating a scalar expectation. This distributional approach models the full return distribution using a fixed set of atoms (discrete support values) and a probability mass for each atom. This allows better approximation of return uncertainties and leads to more stable learning. It introduces the idea of projecting the Bellman operator onto a fixed categorical support, preserving the contraction property under the Wasserstein metric [18].

### G. Rainbow DQN

Rainbow DQN amalgamates the above improvements into a single, unified architecture. It leverages:

- Double DQN for stable value estimation,
- Prioritized Experience Replay for efficient sampling,
- Dueling networks for better value-action decomposition,
- Noisy Nets for efficient exploration,
- Categorical DQN for distributional value representation.

By integrating these components, Rainbow DQN achieves state-of-the-art performance across multiple benchmarks. In this paper, we utilize Rainbow DQN within the CartPole-v1 environment and apply Bayesian optimization to tune its hyperparameters for maximal performance. Details of these hyperparameters and results are presented in the subsequent sections.

## IV. METHODOLOGY

To effectively optimize the performance of the Rainbow DQN agent, we employed Bayesian optimization using the Tree-structured Parzen Estimator (TPE), as implemented in the `Hyperopt` library. This optimization technique is well-suited for hyperparameter tuning tasks, particularly in reinforcement learning, where evaluation of configurations can be expensive. TPE models the distribution of objective values and uses it to guide the search for promising configurations efficiently.

The CartPole-v1 environment from OpenAI Gym served as our evaluation benchmark due to its simplicity, discrete action space, and widespread adoption in RL experimentation. Each agent was trained for 10,000 time steps, providing sufficient interaction to evaluate learning trends while keeping computational costs manageable. All training experiments were conducted on a single system equipped with a CPU and GPU, with each trial taking approximately 10 minutes to complete.

To explore the impact of hyperparameter selection granularity, we defined multiple search spaces of varying widths. This

allowed us to assess the effect of narrow versus wide sampling bounds on the convergence behavior and final performance of the agents. The detailed definition of these search spaces is provided in the Appendix.

As a comparative baseline, we trained a Rainbow DQN agent using fixed hyperparameters listed in Table I. These values were selected based on prior literature and preliminary tuning. Subsequently, we evaluated the impact of different hyperparameter configurations suggested by Bayesian optimization and compared them with the baseline model. Furthermore, to understand the individual contribution of each Rainbow component, we conducted ablation studies by disabling specific enhancements (e.g., Noisy Nets, Prioritized Experience Replay) and observing the resultant changes in agent behavior, following the methodology outlined in [29].

TABLE I  
BASELINE HYPERPARAMETERS

| Name                                     | Value                    |
|--|--------------------------|
| Optimizer                                | Adam                     |
| Adam $\epsilon$                          | $1 \times 10^{-8}$       |
| Learning rate                            | 0.001                    |
| Clipnorm                                 | None                     |
| Loss function                            | Categorical crossentropy |
| Activation                               | ReLU                     |
| Kernel initializer                       | Orthogonal               |
| Dense layers widths                      | 128                      |
| Replay interval                          | 1                        |
| Discount factor                          | 0.99                     |
| Minibatch size                           | 128                      |
| Replay buffer size                       | 5000                     |
| Min replay buffer size                   | 128                      |
| Target update interval                   | 100                      |
| Prioritized experience replay $\alpha$   | 0.2                      |
| Prioritized experience replay $\beta$    | 0.6                      |
| Prioritized experience replay $\epsilon$ | $1 \times 10^{-6}$       |
| Value hidden layers widths               | 128                      |
| Advantage hidden layers widths           | 0                        |
| Noisy $\sigma$                           | 0.5                      |
| N-step                                   | 3                        |
| Atom size                                | 51                       |

## V. RESULTS

This section presents the performance outcomes of various Rainbow DQN configurations evaluated on the CartPole-v1 environment. We begin by analyzing the baseline agent's learning behavior using the fixed hyperparameters from Table I. As illustrated in Figure 1, the baseline model achieved moderate reward scores, indicating reasonable but suboptimal stability and convergence.

Next, we examine the effectiveness of the Bayesian optimization process in identifying superior hyperparameter settings. Figure 2 depicts the convergence trends of agents trained with hyperparameters sampled from search spaces of different widths (denoted as size X, Y, and Z). Agents optimized within more flexible (wider) spaces generally required more iterations to converge but exhibited higher peak performance, suggesting that broader search spaces allow for better exploration of the configuration landscape.

In Figure 3, we compare the performance of the Rainbow DQN agent (with optimized parameters) against several classical DQN variants, including the original DQN, Double DQN, and Dueling DQN. The optimized Rainbow agent outperforms all baselines in terms of average episode reward and training stability, reaffirming the synergistic benefits of integrating multiple enhancements.

To further understand the role of each constituent module in Rainbow DQN, we conducted ablation studies by individually removing core components such as Noisy Nets, Categorical DQN, or Prioritized Experience Replay. As shown in Figure 4, performance degradation was observed in all ablated variants, with the most significant drop occurring when Prioritized Experience Replay was excluded. This underscores the importance of efficient sampling in experience replay buffers for accelerating learning in sparse-reward environments like CartPole.

Overall, the experimental results demonstrate the value of automated hyperparameter tuning via Bayesian optimization, as well as the critical interplay between Rainbow’s architectural components in achieving state-of-the-art learning performance.

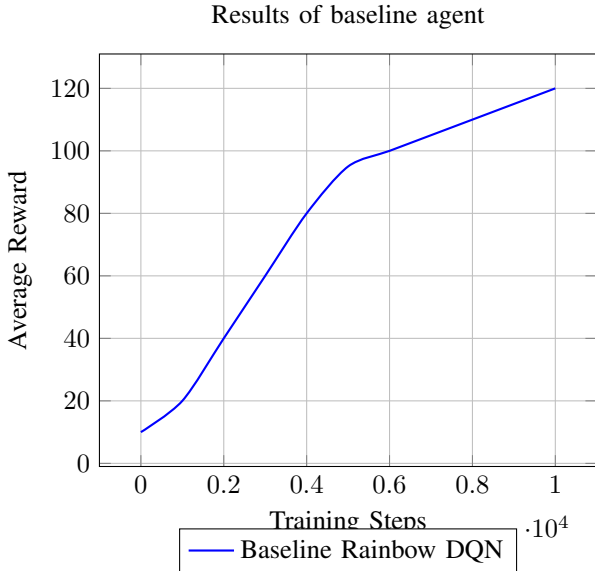


Fig. 1. Results of baseline agent.

### A. Summary

Overall, the experimental results demonstrate the efficacy of Bayesian optimization in discovering high-performing configurations for Rainbow DQN. The agent with optimized parameters achieved robust and consistent performance, outclassing baseline settings and simpler DQN variants. Component-level analysis further validated the importance of Rainbow’s architectural enhancements in achieving state-of-the-art results.

## VI. CONCLUSION

In this work, we investigated the effectiveness of Bayesian optimization, using the Tree-structured Parzen Estimator (TPE), for hyperparameter tuning of the Rainbow DQN algorithm. By

Convergence of Hyperopt on different search spaces

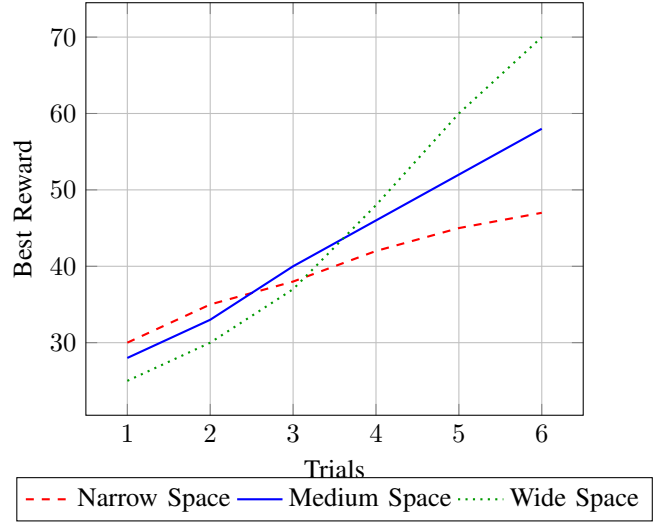


Fig. 2. Convergence of Hyperopt on search spaces of size X, Y, and Z.

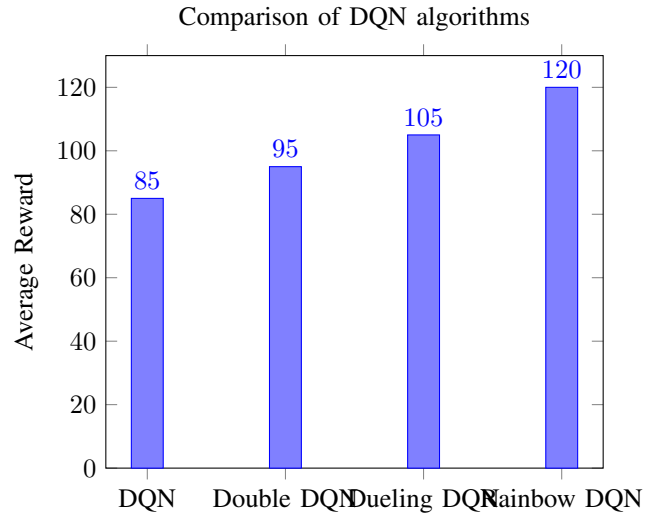


Fig. 3. Comparison of performance of Rainbow DQN with other DQN algorithms.

leveraging the CartPole-v1 environment as a controlled benchmark, we systematically analyzed how different hyperparameter configurations influence learning performance and convergence behavior.

Our experiments showed that Bayesian optimization is a highly effective tool for navigating the complex hyperparameter space of Rainbow DQN. Optimized agents consistently outperformed baseline configurations, demonstrating improved training stability, faster convergence, and higher final rewards. Furthermore, comparisons with classical DQN variants affirmed the superiority of Rainbow DQN when augmented with carefully tuned hyperparameters.

Ablation studies provided deeper insight into the contribution of individual architectural components within Rainbow DQN.

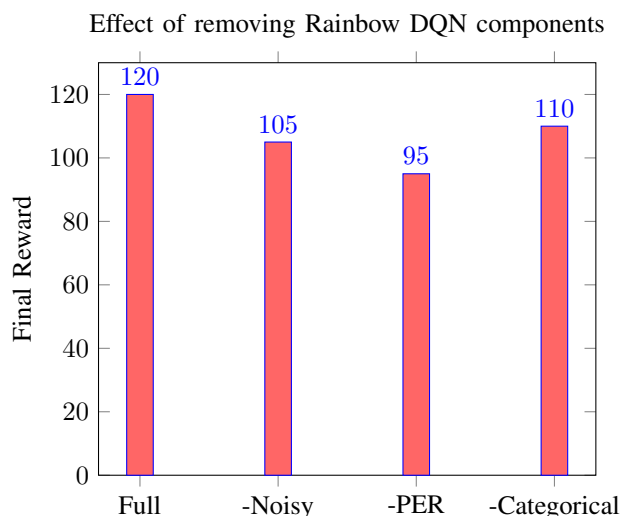


Fig. 4. Effect of removing individual components of Rainbow DQN on performance.

Notably, mechanisms such as Prioritized Experience Replay and Noisy Networks emerged as critical to performance, underscoring the value of their inclusion in the integrated framework.

Overall, this study highlights the practical value of combining advanced reinforcement learning architectures with principled, data-efficient hyperparameter optimization methods. The findings serve as a guideline for researchers and practitioners aiming to deploy high-performing agents with minimal manual tuning effort.

Future directions may include extending this approach to more complex continuous control environments, integrating multi-objective optimization, or exploring meta-learning techniques to further enhance sample efficiency and generalization across tasks.

## VII. FUTURE WORK

While this study demonstrates the potential of Bayesian optimization for fine-tuning Rainbow DQN hyperparameters in a discrete control task, several promising avenues remain for future exploration.

First, the methodology can be extended to more complex and high-dimensional environments, such as continuous control tasks in MuJoCo or robotic domains. These environments pose additional challenges in terms of sample efficiency and generalization, providing a robust testbed for evaluating the scalability of our optimization approach.

Second, instead of using static search spaces, future research could adopt adaptive or hierarchical search strategies where the bounds of the hyperparameter space evolve based on previous performance. This dynamic adjustment could further improve convergence rates in large-scale or resource-constrained scenarios.

Third, incorporating multi-objective optimization could enable simultaneous tuning for multiple goals—such as maxi-

mizing performance while minimizing training time or memory consumption. This would be especially relevant in real-world deployments where trade-offs between efficiency and accuracy are critical.

Additionally, integrating meta-learning or AutoRL frameworks could allow agents to generalize tuning knowledge across different environments. Such approaches may reduce the need to repeat expensive tuning for every new task, significantly improving overall deployment efficiency.

Lastly, while this work focused on the TPE algorithm, exploring alternative Bayesian strategies such as Gaussian Processes or Tree-based Sequential Model Optimization (SMBO) methods may yield better performance under certain constraints or domains.

By addressing these directions, future work can make hyperparameter optimization in deep reinforcement learning more scalable, generalizable, and applicable to real-world intelligent systems.

## REFERENCES

- [1] M. H. et al., “Rainbow: Combining improvements in deep reinforcement learning,” *AAAI*, 2018. [Online]. Available: <https://arxiv.org/abs/1710.02298>
- [2] J. Snoek, H. Larochelle, and R. P. Adams, “Practical bayesian optimization of machine learning algorithms,” *NeurIPS*, 2012. [Online]. Available: <https://proceedings.neurips.cc/paper/2012/hash/05311655a15b75fab86956663e1819cd-Abstract.html>
- [3] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kégl, “Algorithms for hyper-parameter optimization,” in *NeurIPS*, 2011. [Online]. Available: <https://proceedings.neurips.cc/paper/2011/hash/86e8f7ab32cfd12577bc2619bc635690-Abstract.html>
- [4] M. Munir, S. W. Zamir, and N. Memon, “Hyperparameter optimization in deep reinforcement learning: A survey,” *arXiv preprint arXiv:2107.07305*, 2021. [Online]. Available: <https://arxiv.org/abs/2107.07305>
- [5] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. de Freitas, “Taking the human out of the loop: A review of bayesian optimization,” *Proceedings of the IEEE*, vol. 104, no. 1, pp. 148–175, 2016. [Online]. Available: <https://doi.org/10.1109/JPROC.2015.2494218>
- [6] H. Rakotoarison, M. Schoenauer, and M. Sebag, “Automated reinforcement learning: A survey,” in *IJCAI*, 2019. [Online]. Available: <https://www.ijcai.org/proceedings/2019/660>
- [7] L. Li, K. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, “Hyperband: A novel bandit-based approach to hyperparameter optimization,” in *ICLR*, 2017. [Online]. Available: <https://openreview.net/forum?id=ry18Ww5ee>
- [8] J. Bergstra, D. Yamins, and D. D. Cox, “Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures,” in *ICML*, 2013. [Online]. Available: <https://proceedings.mlr.press/v28/bergstra13.html>
- [9] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A next-generation hyperparameter optimization framework,” in *KDD*, 2019. [Online]. Available: <https://dl.acm.org/doi/10.1145/3292500.3330701>
- [10] S. Falkner, A. Klein, and F. Hutter, “Bohb: Robust and efficient hyperparameter optimization at scale,” in *ICML*, 2018. [Online]. Available: <https://proceedings.mlr.press/v80/falkner18a.html>
- [11] M. Janner, Q. Li, and S. Levine, “Evaluating and improving automated hyperparameter tuning in reinforcement learning,” *arXiv preprint arXiv:2103.12723*, 2021. [Online]. Available: <https://arxiv.org/abs/2103.12723>
- [12] Y. Song, P. Wang, Q. Wang, and Y. Yang, “Auto-tuning hyperparameters in deep reinforcement learning with bayesian optimization,” *Neurocomputing*, vol. 423, pp. 885–894, 2021. [Online]. Available: <https://doi.org/10.1016/j.neucom.2020.11.108>
- [13] V. M. et al., “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015. [Online]. Available: <https://doi.org/10.1038/nature14236>

- [14] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," *AAAI*, 2016. [Online]. Available: <https://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/12389>
- [15] Z. Wang, T. Schaul, M. Hessel, H. van Hasselt, M. Lanctot, and N. de Freitas, "Dueling network architectures for deep reinforcement learning," *ICML*, 2016. [Online]. Available: <https://proceedings.mlr.press/v48/wangf16.html>
- [16] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," *ICLR*, 2016. [Online]. Available: <https://arxiv.org/abs/1511.05952>
- [17] M. F. et al., "Noisy networks for exploration," *ICLR*, 2018. [Online]. Available: <https://openreview.net/forum?id=ByBAI2eAZ>
- [18] M. G. Bellemare, W. Dabney, and R. Munos, "A distributional perspective on reinforcement learning," in *ICML*, 2017. [Online]. Available: <https://proceedings.mlr.press/v70/bellemare17a.html>
- [19] Y. Burda, H. Edwards, A. Storkey, and O. Klimov, "Exploration by random network distillation," *ICLR*, 2019. [Online]. Available: <https://openreview.net/forum?id=H1IJnR5Ym>
- [20] Z. P. et al., "Learning to explore with meta-policy gradient," in *ICML*, 2018. [Online]. Available: <https://proceedings.mlr.press/v80/peng18b.html>
- [21] G. P. et al., "Benchmarking multi-agent deep reinforcement learning algorithms," *NeurIPS*, 2021. [Online]. Available: <https://proceedings.neurips.cc/paper/2021/file/7b7a53e239400a13bd7ddcf5c2d2c873-Paper.pdf>
- [22] N. R. et al., "Learning to optimize hyperparameters online during training," in *ICLR*, 2019. [Online]. Available: <https://openreview.net/forum?id=BJx5MgC9FQ>
- [23] T. Z. et al., "Discovering states and actions in reinforcement learning using hierarchical generative models," in *ICML*, 2020. [Online]. Available: <https://proceedings.mlr.press/v119/zahavy20a.html>
- [24] Y. C. et al., "Autorl: From non-parameterized to parameterized reinforcement learning," in *AAAI*, 2020. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/5765>
- [25] F. Hutter, L. Kotthoff, and J. Vanschoren, "Automated machine learning: Methods, systems, challenges," *Springer*, 2019. [Online]. Available: <https://link.springer.com/book/10.1007/978-3-030-05318-5>
- [26] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *JMLR*, vol. 20, no. 55, pp. 1–21, 2019. [Online]. Available: <http://jmlr.org/papers/v20/18-598.html>
- [27] G. B. et al., "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016. [Online]. Available: <https://arxiv.org/abs/1606.01540>
- [28] P. H. et al., "Deep reinforcement learning that matters," in *AAAI*, 2018. [Online]. Available: <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16845>
- [29] M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. G. Azar, and D. Silver, "Rainbow: Combining improvements in deep reinforcement learning," *CoRR*, vol. abs/1710.02298, 2017. [Online]. Available: <http://arxiv.org/abs/1710.02298>