

Reinforcement Learning Quantum Local Search

Chen-Yu Liu ^{*‡}, and Hsi-Sheng Goan [†]

^{*} Graduate Institute of Applied Physics, National Taiwan University, Taipei, Taiwan

[†] Department of Physics, National Taiwan University, Taipei, Taiwan

Email:[‡] d10245003@g.ntu.edu.tw

Abstract—Quantum Local Search (QLS) is a promising approach that employs small-scale quantum computers to tackle large combinatorial optimization problems through local search on quantum hardware, starting from an initial point. However, the random selection of the sub-problem to solve in QLS may not be efficient. In this study, we propose a reinforcement learning (RL) based approach to train an agent for improved sub-problem selection in QLS, beyond random selection. Our results demonstrate that the RL agent effectively enhances the average approximation ratio of QLS on fully-connected random Ising problems, indicating the potential of combining RL techniques with Noisy Intermediate-scale Quantum (NISQ) algorithms. This research opens a promising direction for integrating RL into quantum computing to enhance the performance of optimization tasks.

Index Terms—Reinforcement Learning, Quantum Optimization, Local Search,

I. INTRODUCTION

Exploiting the principles of quantum mechanics, quantum computing facilitates calculations and problem-solving with unparalleled efficiency [1], [2]. This technology holds immense promise for expediting solutions to specific problems beyond the capabilities of classical computers, thereby offering a myriad of application opportunities. To solve problems on quantum computers, it is necessary to map the original problem to a form that can be addressed by quantum algorithms. One common approach for tackling a wide variety of combinatorial optimization problems is to convert them into quadratic unconstrained binary optimization (QUBO) problems, map the problem to an Ising Hamiltonian, and then determine the Hamiltonian’s ground state by variational quantum eigensolver (VQE) [9]–[13] or quantum annealing [3], [5], where the ground state corresponds to the solution of the original problem [3], [4]. Quantum optimization has been successfully applied to real-world challenges, such as portfolio optimization [5]–[8], industrial optimization [14], [15], recommendation system [16], [17], and traveling salesman problems [18], [19].

Given the limitations on the number of qubits during the Noisy Intermediate-scale Quantum (NISQ) era, several methods have been proposed to decrease the number of qubits required for VQE and/or quantum annealing, including cluster-based approaches, large-system sampling approximation, quantum local search (QLS), and other techniques [20]–[26].

The QLS algorithm is designed to address combinatorial optimization problems by iteratively solving sub-problems and

updating local configurations, such that the required number of qubits is just the size of the sub-problem instead of the full problem. This approach enables the use of both gate-based quantum chips and quantum annealers for solving sub-problems. The random selection of sub-problems offers opportunities for enhancement by investigating alternative sub-problem selection strategies, such as the reinforcement learning (RL) strategy explored in this work.

II. PRELIMINARY

In this preliminary section, we review the key ingredients of our work: RL and QLS. RL is a computational approach that enables agents to learn optimal decision-making strategies through interaction with their environment, while QLS is a technique that leverages small quantum computers to solve large combinatorial optimization problems by performing local search on quantum hardware.

A. Quantum Local Search

Quantum Local Search (QLS) [20], [25] combines the principles of quantum computing with local search techniques to solve combinatorial optimization problems. In combinatorial optimization, the goal is to find the optimal solution within a large, discrete solution space, often consisting of various possible configurations. Classical local search [27] algorithms operate by iteratively exploring neighboring solutions in the solution space and moving towards better solutions based on a predefined criterion. QLS utilizes small-scale quantum computers to perform local search operations and leverage the power of quantum computing to enhance the search process. By employing quantum computing, QLS can explore multiple neighboring solutions simultaneously due to quantum parallelism, leading to a potentially more efficient search process.

For an Ising problem with N variables that corresponding to some combinatorial optimization problems:

$$H = \sum_{i,j=1}^N J_{ij} \sigma_i^z \sigma_j^z + \sum_{i=1}^N h_i \sigma_i^z, \quad (1)$$

with J_{ij} the coupling terms and h_i the linear terms. Given an initial starting point of the trial solution state $|\psi\rangle$, QLS iteratively and randomly selects sub-problems with size m such that:

$$H_{\text{sub}} = \sum_{i',j'=1}^m J_{i'j'} \sigma_{i'}^z \sigma_{j'}^z + \sum_{i'=1}^m h_{i'} \sigma_{i'}^z, \quad (2)$$

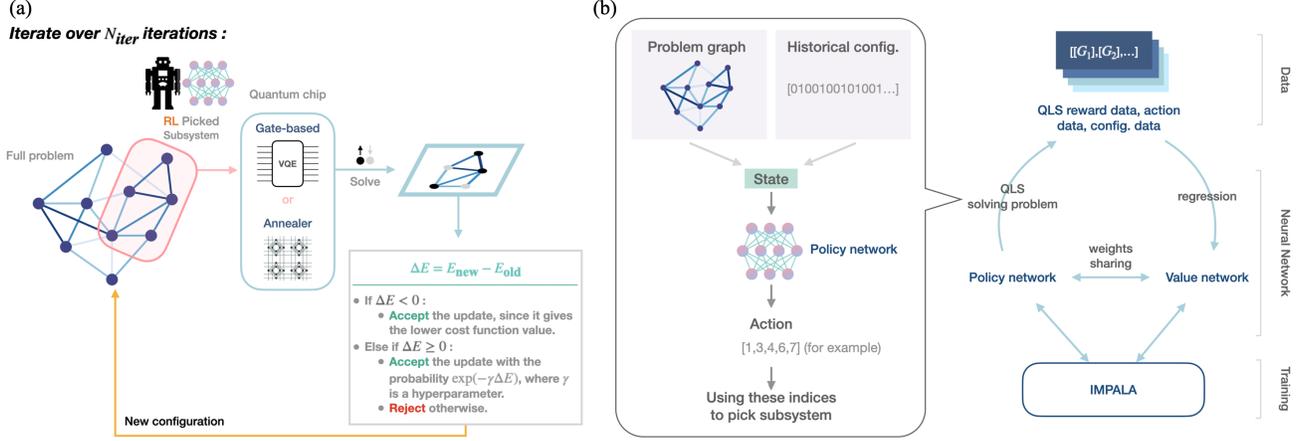


Fig. 1. (a) In the RL-QLS scheme, with N_{iter} iterations, the RL agent selects the sub-problem to be solved by either a gate-based quantum chip or a quantum annealer from the full problem and accepts the trial solution update under specific conditions. (b) The training scheme for RL-QLS involves the policy network/agent determining the sub-problem based on the problem graph information and the historical configuration of trial solutions. It generates data for different problem graphs G_i to train the neural networks using the IMPALA algorithm.

within the larger optimization problem H and solves them using quantum algorithms such as VQE on gate-based quantum computer or quantum annealing on quantum annealer. The algorithm then updates the current solution $|\psi\rangle$ of size N based on the results $|\phi\rangle$ of size m obtained from the quantum computation, moving towards an improved solution. This process is repeated until a satisfactory solution is found or a predefined stopping criterion is met.

The random selection process of subsystems in QLS presents potential inefficiencies in solution exploration. These inefficiencies include redundancy, imbalance, and lack of direction. Redundancy arises when the algorithm chooses similar or already explored sub-problems, wasting computational resources and time. Imbalance occurs when the random selection fails to explore the solution space evenly, possibly overlooking better solutions. The lack of direction results from not utilizing problem-specific knowledge or learning from past iterations to guide the selection process.

To address these issues, alternative approaches, such as reinforcement learning (RL) agents, can be employed to enhance sub-problem selection in QLS. This adaptive approach can improve solution exploration efficiency, enabling the algorithm to discover superior solutions more effectively.

B. Reinforcement Learning

RL [28] is a subfield of machine learning that focuses on training agents to make decisions by interacting with an environment. In RL, an agent learns to perform actions that maximize cumulative rewards over time. The learning process is typically modeled as a Markov Decision Process (MDP), which consists of state space \mathbb{S} , action space \mathbb{A} , reward function, and transition probabilities. The goal is to learn an optimal policy $\pi^*(a|s)$, which dictates the best action $a \in \mathbb{A}$

to take in each state $s \in \mathbb{A}$, that maximizes the expected discounted return:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=0}^T \Gamma^t R_t \right], \quad (3)$$

where R_t is the reward obtained at time step t , and $\Gamma \in [0, 1]$ is a discount factor that determines the relative importance of immediate versus future rewards.

In this work, we use Importance Weighted Actor-Learner Architecture (IMPALA) [29] to train the RL agent. IMPALA is a distributed deep RL algorithm designed to scale up the learning process using multiple parallel actor-learner processes. It achieves this by decoupling the acting and learning processes, allowing for efficient use of computational resources. The actors generate trajectories (s_t, a_t, R_t, s_{t+1}) by interacting with the environment using the current policy, while the learner updates the policy parameters θ using off-policy learning with a value function critic $V(s; \theta_v)$. The policy gradient update is performed using the loss function:

$$\mathcal{L}(\theta) = \mathbb{E}_{\tau} \left[\sum_{t=0}^T \rho_t \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) \delta_t^{\pi_{\theta}} \right], \quad (4)$$

where τ represents a trajectory generated by policy, $\rho_t = \frac{\pi_{\theta}(a_t|s_t)}{\mu_{\theta}(a_t|s_t)}$ is the importance sampling ratio with μ_{θ} the older policy which generates trajectories in replay buffer, and $\delta_t^{\pi_{\theta}} = R_t + \Gamma V(s_{t+1}; \theta_v) - V(s_t; \theta_v)$ is the temporal difference (TD) error.

The use of importance sampling in IMPALA allows for stable and efficient off-policy learning, mitigating the issues of stale trajectories generated by the actors. The distributed nature of the algorithm enables the handling of large-scale environments and high-dimensional state spaces. By combining the advantages of distributed RL with off-policy learning,

IMPALA offers a scalable and powerful solution for a wide range of reinforcement learning problems.

III. REINFORCEMENT LEARNING QUANTUM LOCAL SEARCH

In this work, we propose Reinforcement Learning Quantum Local Search (RL-QLS) to replace the random selection process of sub-problems in QLS with a strategy represented by a trained RL agent, as depicted in the schematic diagram in Fig. 1(a). To train such an agent within the QLS framework, the formulation of fundamental elements of RL, including state, action, reward, and environment, is required. The agent’s state, $s = (G, \vec{\phi})$, is constructed using information from the size- N problem graph $G = (J_{11}, J_{12}, \dots, J_{NN}, h_1, \dots, h_N)$, where J_{ij} represents coupling terms and h_i represents linear terms as in (1), along with the historical configuration $\vec{\phi}$ of the trial solution up to the previous 1 step. The agent’s action $a = (q_1, q_2, \dots, q_m)$ is a list consisting of m elements, where the elements $q_i \in \{0, 1, \dots, N - 1\}$ represent the indices of the variables for the original problem of size N and m is the size of the sub-problem. The reward \mathbf{R} is defined as the approximation ratio $R_{\text{ar}} = \text{TSE}/\text{Dwave-tabu GSE}$, where TSE is the trial solution energy and Dwave-tabu GSE is the ground state energy obtained by the Dwave-tabu solver. The environment dynamics for the RL agent involve a flow from inputting the action to outputting the reward and the subsequent state. In the context of QLS, one step of the environment dynamics entails the following sub-steps:

- 1) For a state $s = (G, \vec{\phi})$, policy $\pi(a|s)$ produces an action a .
- 2) For the sub-problem selected by the action a , compute the sub-problem solution and calculate the corresponding new trial solution energy E_{new} for the full problem.
- 3) With E_{old} representing the old trial solution energy, accept the update if $\Delta E = E_{\text{new}} - E_{\text{old}} < 0$. If $\Delta E \geq 0$, the update is accepted with probability $e^{-\gamma \Delta E}$; otherwise, the update is rejected. Here, γ is a hyperparameter.
- 4) Calculate the reward for the new configuration (trial solution).
- 5) Construct the new state using the above information, $s_{\text{new}} = (G, \vec{\phi}_{\text{new}})$, while the graph information remains unchanged.

A. Training the agent

Following the description of the basic elements of RL in the QLS scheme, the training process for the RL agent using the IMPALA algorithm in this work is set up as follows: For each training iteration, 100 episodes are evaluated, with each episode consisting of 200 steps of quantum local search (environment dynamics) iterations to generate reward data. A set of 1000 randomly generated fully-connected Ising problems serve as the training problem set, such that for each episode, the agent will randomly start from one of these problems. This configuration enables the RL agent to gain a more general understanding of solving various Ising problems. The schematic diagram of the RL training is shown in Fig. 1(b).

The training problem set consists of Ising problems defined as in (1), the range of the couplings and linear terms are $J_{ij} \in (-1, 1)$ and $h_i \in (-1, 1)$. The hyperparameter γ , used to determine the acceptance rate of updates, is set to $\frac{100}{\text{Dwave-tabu GSE}}$, so that the energy difference is normalized for different energy scale. Three cases of size- N Ising problems are examined, where $N \in \{32, 64, 128\}$. The sub-problem size m is set to 5 for all three cases of size- N Ising problems, so that quantum solvers with 5 qubits can be employed. The Dwave-tabu solver is utilized as a quantum simulator for solving the sub-problems to assess the effectiveness of the RL sub-problem selection strategy.

The episode mean reward during the training process is depicted in Fig. 2(a), with the “ N pick m ” labels representing the size- N full problem and corresponding size- m sub-problem. For the “32 pick 5” case, 2000 training iterations are utilized, and the episode mean reward converges to approximately 60 by the end of the training process. In the “64 pick 5” case, 2000 training iterations are also used, exhibiting similar behavior to the “32 pick 5” case, but with an episode mean reward of about 37 at the end. The “128 pick 5” case employs 4300 training iterations, and the episode mean reward approaches 16 by the end of the training process. The observation that cases with larger N values have smaller overall episode mean rewards reflects the increased difficulty associated with solving larger-sized problems.

B. Testing the agent

To evaluate the trained agent, three testing problem sets corresponding to the three cases of problem size N are generated, with couplings and linear terms in the same range as in the training problem set. For the cases “32 pick 5” and “64 pick 5”, $N_t = 100$ test problems are assessed, while for the case “128 pick 5”, $N_t = 30$ test problems are examined. The test results are displayed in Fig. 2(b), with the results of RL-QLS depicted in red and those of the random sub-problem selection strategy shown in blue for comparison. For the cases “32 pick 5” and “64 pick 5”, optimization is performed for 200 steps, and 500 steps are used for the case “128 pick 5”. A noticeable improvement can be observed for the RL-QLS compared to the original random sub-problem selection strategy.

IV. DISCUSSION

The training and testing processes revealed that the RL agent’s performance was affected by the problem size, with smaller episode mean rewards observed for larger problems. This result highlights the increased difficulty of solving larger-sized problems and suggests that further research may be needed to optimize the agent’s performance for these more challenging tasks. The RL-QLS method has shown promise in addressing QLS problems and improving sub-problem selection strategies.

The potential ways to improve the performance of the RL agent could be refining the state representation to incorporate

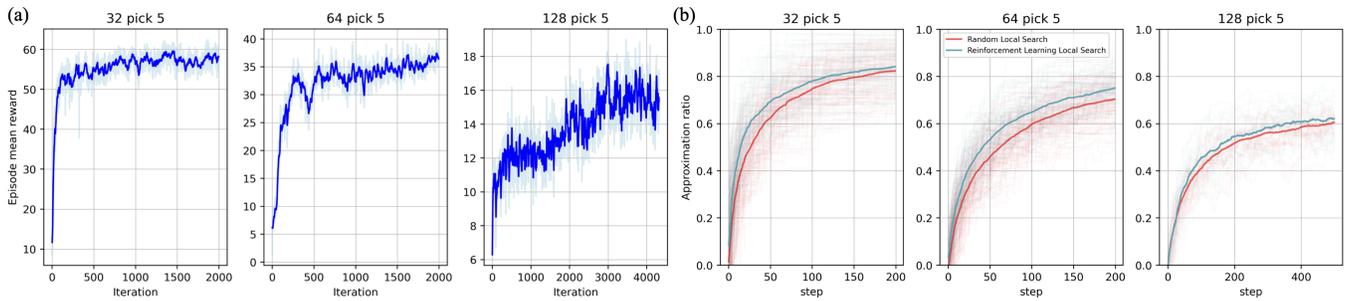


Fig. 2. (a) Episode mean reward during the training process, where “ N pick m ” labels represent the size- N full problem and corresponding size- m sub-problem. (b) The approximation ratio results of RL-QLS (blue) applied to the Ising problems testing set compared with that of random sub-problem selection strategy (red). The solid line represents the average value among the N_t tests, while individual test results are shown as faded lines. For the “32 pick 5” and “64 pick 5” cases, $N_t = 100$, while for the “128 pick 5” case, $N_t = 30$.

more complex problem information or richer historical configurations, and the exploration of alternative RL algorithms and training techniques could potentially enhance the agent’s performance and adaptability, enabling it to tackle more complex and diverse problem sets.

Finally, conducting experiments on real quantum devices would provide a valuable assessment of the practical applicability of the RL-QLS approach in real-world scenarios, highlighting areas for further exploration and enhancement.

V. CONCLUSION

In this work, we introduced RL-QLS, an approach that leverages RL to enhance the sub-problem selection process in QLS. The RL-QLS method replaces the random selection process in QLS with a trained RL agent, which adapts its strategy based on its interaction with the environment. This adaptive approach enables the RL-QLS algorithm to improve the exploration efficiency of combinatorial optimization problems, leading to more effective discovery of superior solutions. We implemented the IMPALA algorithm to train the RL agent, making it capable of handling large-scale environments and high-dimensional state spaces.

The simulation results demonstrated that the RL-QLS algorithm outperformed the random sub-problem selection strategy in solving Ising problems (combinatorial optimization problems), showing the effectiveness of the RL approach in enhancing the QLS method. The episode mean rewards obtained during the training process indicated that the RL agent was able to learn from the environment and improve its performance over time. In the testing phase, the RL-QLS algorithm achieved better approximation ratios compared to the random selection strategy, validating the potential of the proposed method in solving large-scale optimization problems.

In conclusion, the RL-QLS algorithm presents a promising direction for solving problems in combinatorial optimization (Ising problem) by exploiting the capabilities of quantum computing and reinforcement learning. Future work can apply the RL-QLS approach to other optimization problems and examine the scalability of the algorithm as quantum computing resources improve. Integrating more advanced RL techniques

and investigating the transferability of the trained agent can further enhance the performance and efficiency of the RL-QLS method, providing even better solutions to complex optimization challenges.

REFERENCES

- [1] A. Steane, “Quantum computing,” *Reports on Progress in Physics*, vol. 61, no. 2, pp. 117, Feb. 1998.
- [2] J. Preskill, “Quantum Computing in the NISQ Era and Beyond,” *Quantum*, vol. 2, pp. 79, Aug. 2018.
- [3] Z. Bian, F. Chudak, W. Macready, A. Roy, R. Sebastiani, and S. Varotti, “Solving SAT and MaxSAT with a Quantum Annealer: Foundations, Encodings, and Preliminary Results,” *arXiv*, arXiv:1811.02524, Nov. 2018.
- [4] A. Lucas, “Ising formulations of many NP problems,” *Frontiers in Physics*, vol. 2, pp. 5, 2014.
- [5] J. Cohen, A. Khan, and C. Alexander, “Portfolio Optimization of 40 Stocks Using the DWave Quantum Annealer,” *arXiv*, arXiv:2007.01430, Jul. 2020.
- [6] S. Mugel, C. Kuchkovsky, E. S’anchez, S. Fern’andez-Lorenzo, J. Luis-Hita, E. Lizaso, and R. Or’us, “Dynamic portfolio optimization with real datasets using quantum processors and quantum-inspired tensor networks,” *Phys. Rev. Research*, vol. 4, no. 1, pp. 013006, Jan. 2022.
- [7] N. N. Hegade, P. Chandarana, K. Paul, X. Chen, F. Albarr’an-Arriagada, and E. Solano, “Portfolio Optimization with Digitized-Counterdiabatic Quantum Algorithms,” *arXiv*, arXiv:2112.08347, Dec. 2021.
- [8] S. Palmer, S. Sahin, R. Hernandez, S. Mugel, and R. Orus, “Quantum Portfolio Optimization with Investment Bands and Target Volatility,” *arXiv*, arXiv:2106.06735, Jun. 2021.
- [9] A. Peruzzo, J. McClean, P. Shadbolt, et al., “A variational eigenvalue solver on a photonic quantum processor,” *Nature Communications*, vol. 5, no. 1, p. 4213, Jul. 2014.
- [10] A. Kandala, A. Mezzacapo, K. Temme, et al., “Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets,” *Nature*, vol. 549, no. 7671, pp. 242–246, Sep. 2017.
- [11] J. Tilly, et al., “The Variational Quantum Eigensolver: A review of methods and best practices,” *Physics Reports*, vol. 986, pp. 1–128, 2022.
- [12] R. M. Parrish, E. G. Hohenstein, P. L. McMahon, and T. J. Mart’inez, “Quantum Computation of Electronic Transitions Using a Variational Quantum Eigensolver,” *Phys. Rev. Lett.*, vol. 122, no. 23, p. 230401, Jun. 2019.
- [13] S. Mugel, et al., “Dynamic portfolio optimization with real datasets using quantum processors and quantum-inspired tensor networks,” *Phys. Rev. Res.*, vol. 4, no. 1, p. 013006, Jan. 2022.
- [14] S. J. Weinberg, F. Sanches, T. Ide, K. Kamiya, and R. Correll, “Supply Chain Logistics with Quantum and Classical Annealing Algorithms,” *arXiv*, arXiv:2205.04435, May 2022.
- [15] C. Dalyac, L. Henriot, E. Jeandel, W. Lechner, S. Perdrix, M. Porcheron, and M. Veshchezerova, “Qualifying quantum approaches for hard industrial optimization problems. A case study in the field of smart-charging

- of electric vehicles,” *EPJ Quantum Technology*, vol. 8, no. 1, pp. 12, May 2021.
- [16] C.-Y. Liu, H.-Y. Wang, P.-Y. Liao, C.-J. Lai, and M.-H. Hsieh, “Implementation of Trained Factorization Machine Recommendation System on Quantum Annealer,” *arXiv*, arXiv:2210.12953, 2022.
- [17] R. Nembrini, M. Ferrari Dacrema, and P. Cremonesi, “Feature Selection for Recommender Systems with Quantum Computing,” *Entropy*, vol. 23, no. 8, pp. 970, 2021. doi: 10.3390/e23080970.
- [18] A. E. Moylett, N. Linden, and A. Montanaro, “Quantum speedup of the traveling-salesman problem for bounded-degree graphs,” *Phys. Rev. A*, vol. 95, no. 3, pp. 032323, Mar 2017.
- [19] Ö. Salehi, A. Glos, and J. A. Mischczak, “Unconstrained binary models of the travelling salesman problem variants for quantum optimization,” *Quantum Information Processing*, vol. 21, no. 2, pp. 67, Jan 2022.
- [20] R. Shaydulin, H. Ushijima-Mwesigwa, I. Safro, S. Mniszewski, and Y. Alexeev, “Network Community Detection on Small Quantum Computers,” *Advanced Quantum Technologies*, vol. 2, no. 9, p. 1900029, 2019.
- [21] J. Li, M. Alam and S. Ghosh, “Large-scale Quantum Approximate Optimization via Divide-and-Conquer,” *arXiv*, arXiv:2102.13288, 2021.
- [22] Y. Zhang *et al.*, “Variational Quantum Eigensolver with Reduced Circuit Complexity,” *arXiv*, arXiv:2106.07619, 2021.
- [23] Z. Zhou *et al.*, “QAOA-in-QAOA: solving large-scale MaxCut problems on small quantum machines,” *arXiv*, arXiv:2205.11762, 2022.
- [24] B. Tan *et al.*, “Qubit-efficient encoding schemes for binary optimisation problems,” *Quantum*, vol. 5, p. 454, May 2021.
- [25] M. Booth and S. P. Reinhardt, “Partitioning optimization problems for hybrid classical/quantum execution,” *D-Wave Technical Report Series*, vol. 14-1006A-A, 2017.
- [26] C.-Y. Liu and H.-S. Goan, “Hybrid Gate-Based and Annealing Quantum Computing for Large-Size Ising Problems,” *arXiv*, arXiv:2208.03283, 2022.
- [27] E. Aarts and J. K. Lenstra, Eds., *Local Search in Combinatorial Optimization*, Princeton University Press, Princeton, 2003.
- [28] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, “Deep Reinforcement Learning: A Brief Survey,” *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26-38, 2017.
- [29] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley, I. Dunning, S. Legg, and K. Kavukcuoglu, “IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures,” *arXiv*, arXiv:1802.01561, 2018.