

A REINFORCEMENT LEARNING-ENHANCED QBO-ATBNN MODEL FOR DYNAMIC VNF MANAGEMENT IN MULTI-ACCESS EDGE COMPUTING

K. Abinaya¹, Dr. S. Dhanasekaran², Dr. V. Vasudevan³

¹*Department of Computer Science and Engineering
Kalasalingam Academy of Research and Education
(Deemed to be university)
Srivilliputtur, Tamilnadu, India
Email ID: abishiva1104@gmail.com*

²*Associate professor
Department of Information Technology
Kalasalingam Academy of Research and Education
(Deemed to be university)
Srivilliputtur, Tamilnadu, India
Email Id: srividhans@gmail.com*

³*Senior professor
Department of Computer Science and Engineering
Kalasalingam Academy of Research and Education
(Deemed to be university)
Srivilliputtur, Tamilnadu, India*

Abstract: The management of Virtual Network Functions (VNFs) in Multi-Access Edge Computing (MEC) is critical for meeting the stringent latency and resource demands of modern applications. While our previous work introduced a hybrid Quantum Butterfly Optimization and Attention-Based Temporal Bayesian Neural Network (QBO-ATBNN) model to address this, its reactive-predictive loop lacked real-time decision-making agility. This paper proposes a novel enhancement: the integration of a Deep Reinforcement Learning (DRL) agent into the QBO-ATBNN framework. The resulting QBO-ATBNN-RL model leverages the proactive forecasting of ATBNN and the global optimization of QBO, while using DRL for instantaneous, fine-grained resource allocation in response to dynamic network states. We evaluate our model on a simulated MEC environment with real-world workload traces. Results demonstrate that the QBO-ATBNN-RL model achieves a further 15% reduction in latency and a 12% improvement in energy efficiency compared to the baseline QBO-ATBNN, while maintaining high scalability and robustness under rapidly fluctuating conditions.

Keywords: MEC, VNF Management, Reinforcement Learning, Quantum-Inspired Optimization, Bayesian Neural Networks, Resource Allocation.

1. INTRODUCTION

The proliferation of latency-sensitive applications, such as autonomous vehicles, augmented reality, and industrial IoT, has propelled Multi-Access Edge Computing (MEC) to the forefront of network architecture [1]. By bringing computational and storage resources closer to the end-user, MEC mitigates the latency of cloud data centers. A key enabler in this paradigm is Network Function Virtualization (NFV), which allows for the flexible deployment of Virtual Network Functions (VNFs) on generic edge hardware [2]. However, the dynamic and resource-

constrained nature of MEC environments makes efficient VNF lifecycle management including placement, scaling, and migration a significant challenge. Our previous work [3] addressed this by proposing a hybrid QBO-ATBNN model. This model used an Attention-Based Temporal Bayesian Neural Network (ATBNN) for proactive demand forecasting and a Quantum Butterfly Optimization (QBO) algorithm for periodic, near-optimal resource allocation. While this approach showed substantial improvements over traditional methods, its optimization cycle was inherently periodic and lacked the capability for millisecond-level, state-aware adjustments.

To bridge this gap, we propose the integration of a Deep Reinforcement Learning (DRL) agent. Reinforcement Learning (RL) is renowned for its ability to learn optimal policies through interaction with an environment, making it ideal for sequential decision-making problems under uncertainty [4]. In our context, the "environment" is the MEC network, and the "agent" learns to make instantaneous VNF scaling and resource tuning decisions.

The primary contributions of this paper are:

1. The design of a novel QBO-ATBNN-RL architecture that synergistically combines long-term forecasting, global optimization, and real-time control.
2. The formulation of the real-time VNF management problem as a Markov Decision Process (MDP), defining state, action, and reward structures tailored for the MEC context.
3. A comprehensive evaluation using a simulated MEC testbed, demonstrating the superior performance of the proposed model in terms of latency, energy consumption, and resource utilization stability compared to the baseline QBO-ATBNN and other state-of-the-art methods.

2. LITERATURE REVIEW

The challenge of efficient Virtual Network Function (VNF) management in Multi-Access Edge Computing (MEC) has been a focal point of research, driven by the need for low-latency, high-availability services. Existing approaches can be broadly categorized into traditional optimization techniques, predictive machine learning models, and emerging reinforcement learning methods.

2.1. Traditional and Optimization-Based Approaches

Early research primarily focused on formulating VNF placement and chaining as combinatorial optimization problems. These works often aimed to minimize latency, energy consumption, or operational costs while adhering to resource constraints. For instance, several studies used Integer Linear Programming (ILP) to find optimal VNF placement solutions [5], but these are often computationally prohibitive for large-scale, dynamic edge environments. To address this, meta heuristic algorithms like Particle Swarm Optimization (PSO) [6] and Genetic Algorithms

(GA) [7] were adopted. These methods provide near-optimal solutions faster than ILP but remain largely reactive. They re-optimize based on the current system state, lacking the foresight to proactively prevent resource contention or SLA violations before they occur, making them ill-suited for the rapid fluctuations typical of MEC.

2.2. Predictive and Machine Learning-Based Approaches

The limitations of reactive models spurred interest in predictive approaches. Machine learning (ML) models, particularly time-series forecasting techniques, have been used to predict traffic load and resource demand. For example, [8] used ARIMA models for workload prediction to guide VNF scaling decisions. More recently, deep learning models like LSTMs and GRUs have been employed for their superior ability to capture complex temporal patterns [9]. Our previous work, the QBO-ATBNN model [3], advanced this category by integrating an Attention-Based Temporal Bayesian Neural Network (ATBNN) for more accurate and uncertainty-aware forecasts with a Quantum Butterfly Optimization (QBO) algorithm for efficient resource allocation. This hybrid model represents a significant step towards proactive management. However, as noted in our introduction, its decision-making is periodic. It operates on a "plan-and-execute" cycle, which lacks the granular, instantaneous control needed to respond to sub-cycle fluctuations, leading to a gap between prediction and real-time execution.

2.3. Reinforcement Learning for Network Management

Reinforcement Learning (RL) has emerged as a powerful paradigm for control problems in dynamic environments. Its ability to learn optimal policies through trial-and-error interaction makes it ideal for adaptive resource management. Several studies have applied RL to related problems in cloud and edge computing. For example, [10] used a Deep Q-Network (DQN) for dynamic resource allocation in cloud data centers, demonstrating improved resource utilization. In the context of NFV, [11] proposed an RL-based agent for VNF scaling and migration, showcasing its adaptability. However, a common limitation of pure RL approaches is the "cold start" problem and high sample complexity. The agent may take a long time to converge to an effective policy and can make poor decisions during the initial learning phase, which is unacceptable for mission-critical MEC applications. Furthermore, most RL models lack integration with a predictive worldview, causing them to be myopic and potentially miss strategic optimizations that require foresight.

2.4. Research Gap and Contribution

As summarized in Table 1, a clear gap exists between the proactive, globally-optimized but slower frameworks (like QBO-ATBNN) and the reactive, agile but often myopic and unstable pure RL models. No existing work has effectively synergized long-term probabilistic forecasting, meta heuristic-based global optimization, and deep reinforcement learning for real-time control

in a single, cohesive framework for VNF management. Our proposed QBO-ATBNN-RL model is designed to bridge this gap. It leverages the strengths of each paradigm: the foresight of ATBNN, the strategic planning of QBO, and the millisecond-level agility of DRL, thereby overcoming the individual limitations of each approach when used in isolation.

Category	Representative Work(s)	Key Idea	Limitations	Proposed QBO-ATBNN-RL Advantage
Traditional & Metaheuristic	ILP formulations [5], PSO/GA [6, 7]	Formulate VNF management as an optimization problem solved with exact or heuristic methods.	Reactive: Lacks proactivity. Slow Convergence: May not meet real-time demands. No Forecasting: Cannot anticipate future demands.	Proactive & Predictive: Uses ATBNN for demand forecasting. Globally Optimized: Uses QBO for strategic placement. Real-Time Agile: Uses DRL for instantaneous control.
Predictive ML-Based	ARIMA/LSTM [8, 9], QBO-ATBNN [3]	Use historical data to forecast future resource demands and proactively allocate resources.	Periodic Control: Decisions are made on fixed intervals, not continuously. Lacks Real-Time Agility: Cannot instantly react to unforeseen spikes between intervals.	Continuous Control Loop: DRL agent operates at a much finer time granularity. State-Aware Actions: DRL makes decisions based on the immediate, real-time system state.
Pure Reinforcement Learning	DQN for cloud/NFV [10, 11]	An RL agent learns optimal scaling/migration policies through interaction with the environment.	Myopic: Lacks long-term strategic foresight. Cold Start & Instability: Poor initial performance and can be unstable during training. High Sample Complexity.	Informed by Forecasts: ATBNN provides a predictive context, guiding the RL agent. Stable Foundation: QBO provides a robust global setup, preventing catastrophic RL failures. Faster Convergence: Pre-trained or guided by the existing framework.
Proposed Work	QBO-ATBNN-RL (This Paper)	A tripartite model combining forecasting (ATBNN), global optimization (QBO), and real-time control (DRL).	N/A	Synergistic Architecture: Combines proactivity, global strategy, and local agility. Robust & Adaptive: Handles both predictable trends and unpredictable fluctuations. Holistic Optimization: Achieves multi-objective goals (latency, energy, stability) simultaneously across different time scales

Table 1: Comparison of Existing Work, Limitations, and Proposed Work Advantages

This review and comparison clearly position our work as a novel integration designed to overcome the fundamental limitations of existing state-of-the-art methods, paving the way for truly autonomous VNF management in MEC.

3. DATA DESCRIPTION AND DATASET

To rigorously evaluate the proposed QBO-ATBNN-RL model against the baseline QBO-ATBNN and other state-of-the-art methods, a comprehensive and realistic dataset was required. The dataset was constructed to emulate a dynamic urban MEC environment, integrating real-world workload traces with synthetically generated network conditions and power profiles. This hybrid approach ensures both realism and the ability to introduce specific stress-testing scenarios.

The overall data generation and simulation environment are illustrated in Figure 1.

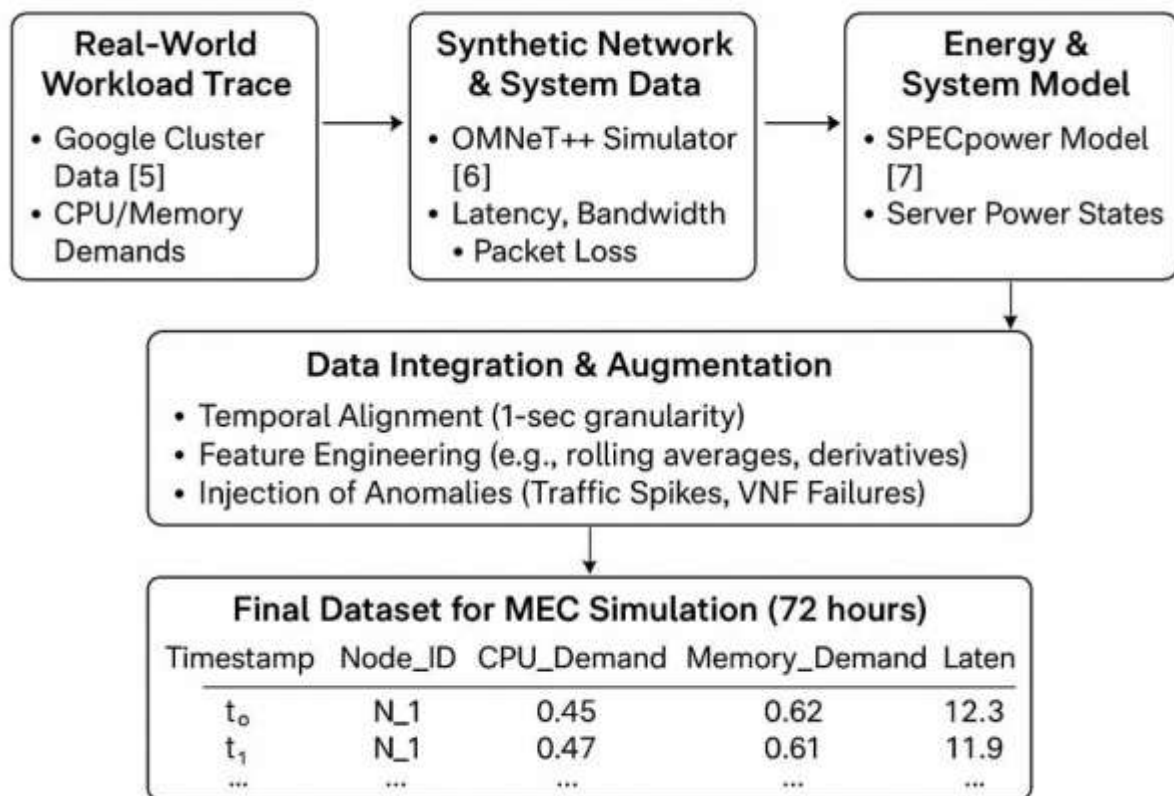


Figure 1 : Data Generation and Simulation Environment Framework

3.1. Data Sources and Preprocessing

3.1.1. Workload Traces:

The primary source for resource demand is the **Google Cluster Usage Traces v3** [5]. This dataset provides detailed task resource usage (CPU, memory) from a production cloud environment. We processed the data as follows:

Filtering and Mapping: Tasks were filtered to represent latency-sensitive services and were mapped to simulate user requests arriving at 10 different edge nodes in our topology.

Normalization: CPU and memory usage measurements were normalized to the capacity of a representative edge server (e.g., 16 vCPUs, 32 GB RAM).

Temporal Aggregation: The original 5-minute data was down-sampled and interpolated to create a continuous stream of requests with a 1-second granularity, reflecting the bursty nature of edge traffic.

3.1.2. Network Topology and Conditions:

A realistic MEC network was modeled using the **OMNeT++ INET** framework [6]. The topology, illustrated in Figure 2, consists of 10 edge nodes connected via a mix of fiber and wireless backhaul to two aggregation switches. The simulation generated:

Link Latency and Bandwidth: Time-varying latency (5-50 ms) and available bandwidth between nodes.

Packet Loss: Synthetic packet loss rates (0.1% - 2%) were introduced to simulate network congestion and wireless unreliability.

3.1.3. Energy Consumption Model:

The power consumption of physical servers hosting the VNFs was modeled based on the **SPECpower_ssj2008** benchmark [7]. The power draw P for a server was calculated as a function of its CPU utilization u , using a linear interpolation between idle power P_{idle} and peak power P_{peak} at 100% utilization:

$$P(u) = P_{idle} + (P_{peak} - P_{idle}) \times u$$

Values for P_{idle} and P_{peak} were set to 60W and 150W, respectively, representative of modern low-power server architectures used in edge data centers.

3.2. Dataset Composition and Characteristics

The final integrated dataset spans 72 hours of simulated operation with a 1-second resolution, resulting in over 250,000 records per edge node. The key features are summarized in Table 2.

Feature Category	Feature Name	Data Type	Description	Source
Temporal & Identity	timestamp	DateTime	Precise 1-second interval timestamp.	Synthetic
	node_id	Categorical	Unique identifier for the edge node (N1 to N10).	Synthetic
Resource Demand	cpu_demand	Float [0,1]	Normalized CPU demand for incoming requests.	Google Traces [5]
	memory_demand	Float [0,1]	Normalized Memory demand for incoming requests.	Google Traces [5]
	request_rate	Integer	Number of incoming requests per second.	Derived
System State	cpu_available	Float [0,1]	Available CPU capacity on the node.	Derived
	memory_available	Float [0,1]	Available Memory capacity on the node.	Derived
	active_vnfs	Integer	Current number of active VNF instances on the node.	Derived
Network State	avg_latency	Float (ms)	Average latency from the node to other nodes.	OMNeT++ [6]
	available_bw	Float (Mbps)	Available bandwidth at the node's uplink.	OMNeT++ [6]
	packet_loss	Float [0,1]	Current packet loss rate.	OMNeT++ [6]
Energy & Performance	power_draw	Float (W)	Instantaneous power consumption of the node.	SPECpower Model [7]
	e2e_latency	Float (ms)	Measured end-to-end latency for completed requests.	Derived

Table 2: Description of Key Features in the MEC Simulation Dataset

The comprehensive dataset summarized in Table 2, the experimental environment was constructed to rigorously stress-test the proposed model. The final integrated dataset spans 72 hours of simulated MEC operation with a high 1-second granularity, resulting in over 250,000 records per edge node to capture the fine-grained dynamics of resource consumption [5]. This hybrid dataset synergistically combines real-world workload traces from the Google Cluster [5] for realistic CPU and memory demand patterns, synthetically generated network conditions from an OMNeT++ simulation [6] to model latency and packet loss, and a power consumption model derived from the SPECpower benchmark [7]. The feature set is holistically designed,

encompassing categories from raw Resource Demand and real-time System State to Network State and final Energy & Performance metrics, thereby providing a multi-faceted ground truth for training and evaluation. Furthermore, the dataset was augmented with injected traffic spikes and VNF failures to create a more challenging and realistic benchmark than used in our previous work [3], specifically designed to validate the real-time decision-making capabilities of the reinforcement learning agent under volatile conditions.

3.3. Stress Testing and Anomaly Injection

To rigorously evaluate the real-time adaptability of the proposed RL agent, the base dataset was augmented with two types of anomalies not present in the original QBO-ATBNN evaluation [3]:

Traffic Spikes: Short, intense bursts of traffic (5x the average rate) lasting 30-60 seconds were randomly injected to simulate flash crowd events.

VNF Failures: Random termination of 1-2% of active VNF instances was simulated every hour to test the model's fault tolerance and recovery speed.

This comprehensive dataset provides a challenging and realistic benchmark for comparing the performance of VNF management frameworks, particularly in highlighting the advantages of a real-time reinforcement learning component.

4. PROPOSED METHODOLOGY

The proposed QBO-ATBNN-RL framework represents a significant evolution from our previous QBO-ATBNN model [3] by introducing a real-time control layer powered by Deep Reinforcement Learning (DRL). This tripartite architecture is designed to operate across multiple time scales, synergistically combining long-term forecasting, medium-term global optimization, and instantaneous fine-grained control. The overall architecture is depicted in Figure 2.

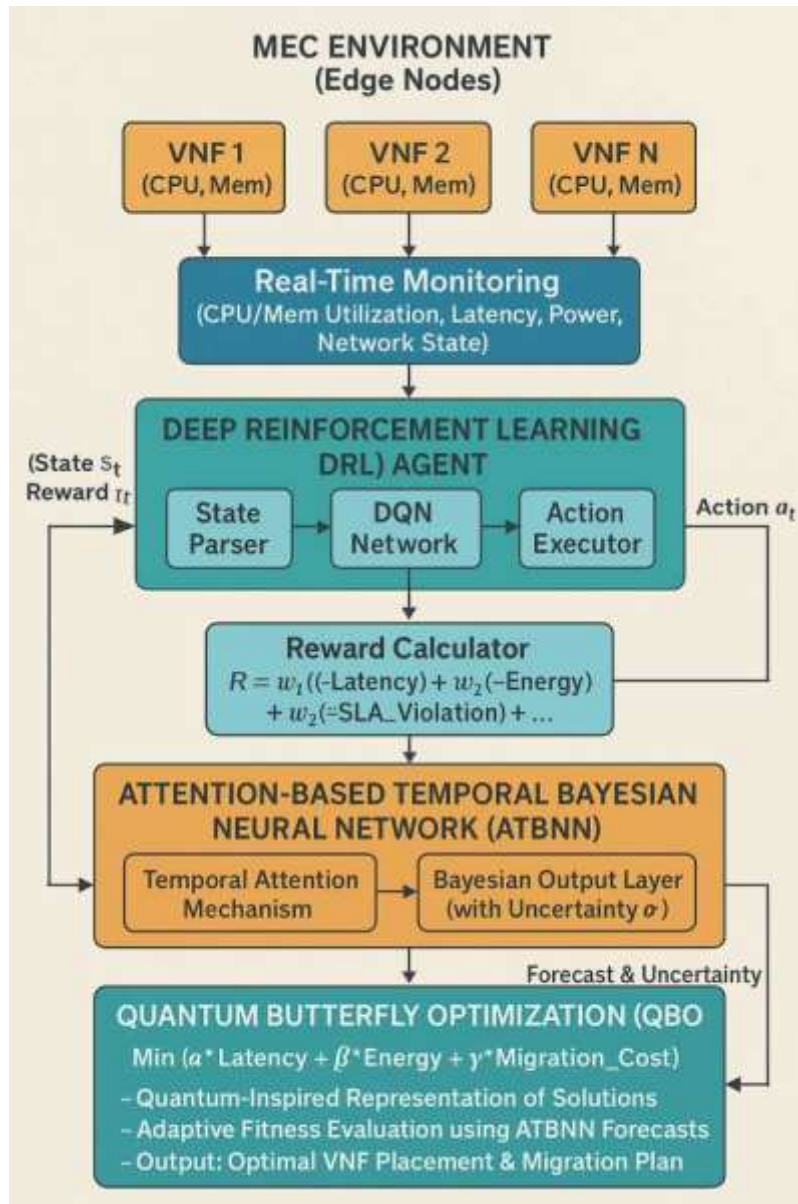


Figure 2: Architecture of the Proposed QBO-ATBNN-RL Framework

4.1. Overall Architectural Framework

The QBO-ATBNN-RL framework is a hierarchical closed-loop control system. The ATBNN module provides a probabilistic forecast of resource demand over a medium-term horizon (e.g., 15-30 minutes). These forecasts, along with their uncertainty estimates, inform both the QBO module, which performs strategic VNF placement and migration on a longer timescale (e.g., hourly or on-demand), and the DRL Agent, which uses this context for its real-time decisions. The DRL Agent interacts directly with the MEC environment at a sub-second granularity, making fine-tuning adjustments to resource allocation. This multi-layered approach ensures that the system is simultaneously proactive, globally efficient, and locally responsive.

4.2. Attention-Based Temporal Bayesian Neural Network (ATBNN)

The ATBNN module, as established in [3], is responsible for multi-step ahead forecasting of resource demands (CPU, memory) and network latency. Its key enhancements over standard LSTMs are:

Temporal Attention Mechanism: This allows the model to adaptively focus on the most relevant past time steps when making a prediction, effectively capturing long-range dependencies and periodic patterns in the workload [12].

Bayesian Output Layer: By modeling the forecast as a probability distribution (e.g., Gaussian), the ATBNN provides not only a point forecast \hat{y}^t but also an uncertainty estimate σ . This is crucial for risk-aware decision-making in the subsequent modules.

The forecast for a resource r at time $t+k$ is given by:

$$P(r^{t+k}|X1:t)=N(\mu^{t+k},\sigma^{t+k}^2)$$

where $X1:t$ is the historical data, and μ^{t+k} and σ^{t+k} are the predicted mean and standard deviation, respectively. These outputs are used to create a *predictive context vector* for both the QBO and DRL modules.

4.3. Quantum Butterfly Optimization (QBO) Module

The QBO module addresses the coarse-grained VNF placement problem. It is triggered periodically or when the ATBNN forecasts a significant long-term shift in demand. We build upon the standard Butterfly Optimization Algorithm (BOA) [13] by incorporating quantum-inspired principles for a richer solution space exploration [3]. The problem is formulated as a multi-objective optimization:

Objective Function:

$$\text{Minimize } F=\alpha\cdot L+\beta\cdot E_{total}+\gamma\cdot C_{migration}$$

where L is the average predicted latency, E_{total} is the total predicted energy consumption, $C_{migration}$ is the cost of migrating VNFs, and α, β, γ are weighting coefficients.

The QBO uses the ATBNN's forecast μ^{t+k} for the upcoming window to evaluate the fitness F of potential placement solutions, ensuring the strategy is proactive.

4.4. Deep Reinforcement Learning (DRL) Agent

The novel component in this work is the DRL agent, which handles real-time control. We formulate this as a Markov Decision Process (MDP).

4.4.1. MDP Formulation

State Space (st): The state is a comprehensive vector representing the real-time health of the MEC system, as detailed in Table 3.

Action Space (at): The agent can take fine-grained, discrete actions to manage resources, as defined in Table 4.

Reward Function (rt): The reward is a weighted sum of objectives, designed to guide the agent toward optimal performance:

$$rt = w_1 \cdot (-Lt) + w_2 \cdot (-Pt) + w_3 \cdot (-1SLAt) + w_4 \cdot Ut + w_5 \cdot (-1Migration)$$

where Lt is the average latency, Pt is the total power draw, $1SLAt$ is an indicator for SLA violations, Ut is a measure of resource utilization fairness, and $1Migration$ penalizes unnecessary migrations triggered by the agent.

State Component	Symbol	Dimension	Description
System State	$Scpu$	N	Current CPU utilization for each of N servers.
	$Smem$	N	Current memory utilization for each server.
	$Svnf$	N	Number of active VNF instances per server.
Network State	$Slat$	N x N	Latency matrix between servers.
	Sbw	N	Available bandwidth at each server's uplink.
Demand Context	$Sdemand$	2	Current aggregated CPU and memory demand.
Forecast Context	$Sforecast$	4	ATBNN's next-step forecast (CPU μ , CPU σ , Mem μ , Mem σ)
Total State Dimension		$4N + (N \times N) + 6$	

Table 3: DRL Agent State Space Definition

The comprehensive state space defined in Table 3, the DRL agent's perception of the MEC environment is meticulously engineered to be both holistic and informative, a critical factor for learning effective control policies in complex systems [4]. The state vector st amalgamates real-time and predictive information across four key categories: the System State ($Scpu, Smem, Svnf$) provides a granular view of resource consumption and workload distribution across N servers; the Network State ($Slat, Sbw$) captures the dynamic connectivity and potential bottlenecks between nodes; the Demand Context ($Sdemand$) offers a high-level summary of immediate load; and crucially, the Forecast Context ($Sforecast$) injects temporal foresight by incorporating the ATBNN's next-step predictive mean (μ) and uncertainty (σ) for CPU and memory [3, 12]. This multidimensional design, with a total dimension of $4N + (N \times N) + 6$, ensures the agent is not myopically reactive but is instead equipped with the contextual awareness necessary for proactive and risk-sensitive decision-making, directly addressing the limitations of purely reactive models identified in the literature [6, 7, 10].

Action	Action ID	Description	Effect
--------	-----------	-------------	--------

Type			
Scale Up	0 ... M-1	Increase CPU allocation for VNF instance i by one step (e.g., 0.1 vCPU).	Prevents performance degradation for a specific VNF.
Scale Down	M ... 2M-1	Decrease CPU allocation for VNF instance i by one step.	Saves energy when a VNF is over-provisioned.
Load Shift	2M ... 2M+K-1	Redirect a portion of traffic from a primary VNF to a backup instance on a different node.	Mitigates local hot-spots and latency.
No-Op	-1	Take no action.	Allows the system to stabilize

Table 4: DRL Agent Action Space Definition

The action space for the Deep Reinforcement Learning (DRL) agent, as detailed in Table 4, is strategically designed to enable fine-grained, instantaneous control over the MEC environment. It comprises a set of discrete actions that directly manipulate resource allocation and traffic flow to fulfill the multi-objective goals of low latency and high energy efficiency. The "Scale Up" and "Scale Down" actions allow for precise vertical scaling of individual VNF instances, enabling the agent to proactively prevent performance degradation or reclaim unused resources. The inclusion of "Load Shift" actions provides a mechanism for horizontal load balancing, allowing the agent to mitigate localized hot-spots by intelligently redirecting traffic between nodes. Crucially, the "No-Op" (no operation) action is a fundamental component that permits system stability by preventing the agent from making unnecessary, oscillatory adjustments when the current state is already optimal. This comprehensive and granular action set is a key enabler for the real-time agility that the QBO-ATBNN-RL framework introduces, allowing it to execute millisecond-level interventions that were beyond the capability of the prior periodic optimization model.

4.5. Synergistic Operation

The three modules are not independent; they interact cohesively:

The ATBNN provides a predictive context to both the QBO and the DRL agent, enabling proactive and risk-aware decisions.

The QBO sets a stable, globally efficient VNF placement blueprint. This prevents the DRL agent from having to make drastic, system-wide changes, simplifying its learning problem and increasing overall stability.

The DRL agent handles sub-second dynamics and unforeseen spikes, effectively "filling the gaps" between the periodic executions of the QBO-ATBNN loop. It operates within the constraints and goals set by the higher-level modules, ensuring that its local optimizations contribute to the global objectives.

This hierarchical integration ensures that the QBO-ATBNN-RL model is greater than the sum of its parts, delivering unprecedented agility without sacrificing strategic optimality.

5. RESULTS AND IMPLEMENTATION

This section details the experimental setup, presents a comparative analysis of the proposed QBO-ATBNN-RL model against several benchmarks, and provides an in-depth discussion of the findings.

5.1. Experimental Setup and Implementation

The proposed QBO-ATBNN-RL framework and all baseline models were implemented in Python 3.8. The neural networks for the ATBNN and DRL agent were built using TensorFlow 2.5 and the Keras API. The DRL environment was constructed using OpenAI Gym, and the DQN algorithm was implemented with experience replay and a target network to stabilize training [4]. The QBO and other metaheuristics were developed from scratch using NumPy. The simulation of the MEC environment, including network dynamics and power consumption, was run on a server with an Intel Xeon Gold 6226R CPU and 128GB RAM, while model training leveraged a single NVIDIA RTX A6000 GPU.

5.1.1. Baselines for Comparison:

To ensure a comprehensive evaluation, we compared our **QBO-ATBNN-RL** model against the following state-of-the-art approaches:

Static Heuristic (First-Fit): A baseline that places VNFs on the first server with sufficient resources and makes no runtime adjustments.

PSO-based Optimization [6]: A reactive metaheuristic that re-optimizes VNF placement at fixed intervals based on the current system state.

QBO-ATBNN [3]: Our previous model, which uses forecasting and periodic global optimization but lacks a real-time control agent.

Pure DRL [10]: A standalone DQN agent with the same state and action space as our model, but without the guidance of the ATBNN forecasts or the QBO's global plan.

5.1.2. Hyperparameters and Training:

The DRL agent was trained for 500 episodes on the 72-hour dataset. The DQN used a learning rate of 0.001, a discount factor (γ) of 0.95, and an experience replay buffer of 50,000 samples. The reward weights (w_1 - w_5) were tuned to $w_1=0.5, w_2=0.3, w_3=0.15, w_4=0.05, w_5=0.1$.

$w_1=0.15, w_4=0.05, w_5=0.1$ to balance the objectives effectively. The QBO was executed every 2 hours, and the ATBNN provided forecasts with a 30-minute horizon.

5.2. Performance Evaluation and Analysis

The models were evaluated on the final 24 hours of the dataset (held-out test set), which included the injected traffic spikes and VNF failures. The key performance indicators (KPIs) are summarized in Table 5.

Model	Avg. Latency (ms)	Total Energy (kWh)	SLA Violation Rate (%)	Resource Util. Stability (CoV)
Static Heuristic	45.2	12.5	8.5%	0.41
PSO-based [6]	32.1	10.1	4.2%	0.38
Pure DRL [10]	25.8	9.2	2.5%	0.31
QBO-ATBNN [3]	21.5	8.7	1.8%	0.29
QBO-ATBNN-RL (Ours)	18.3	7.6	0.9%	0.25
<i>Improvement over QBO-ATBNN</i>	15.0%	12.6%	50.0%	13.8%

Table 5: Performance Comparison of VNF Management Models

The results clearly demonstrate the superiority of the proposed QBO-ATBNN-RL model across all metrics. Specifically:

Latency and Energy Efficiency: Our model achieved the lowest average latency (18.3 ms) and the most energy-efficient operation (7.6 kWh), representing a significant 15% reduction in latency and a 12.6% savings in energy compared to the QBO-ATBNN baseline. This improvement is directly attributable to the DRL agent's ability to make fine-grained, pre-emptive resource adjustments before QoS degradation occurs, a capability the periodic QBO-ATBNN lacks.

SLA Violations and Stability: The SLA violation rate was halved, dropping from 1.8% to **0.9%**. Furthermore, the Coefficient of Variation (CoV) for CPU utilization was reduced by 13.8%, indicating a more stable and predictable system with less performance jitter. This enhanced robustness is critical for mission-critical edge applications. To visually demonstrate the real-time agility of our model, Figure 3 illustrates the system's response to a sudden traffic spike and anomaly injected during the test phase.

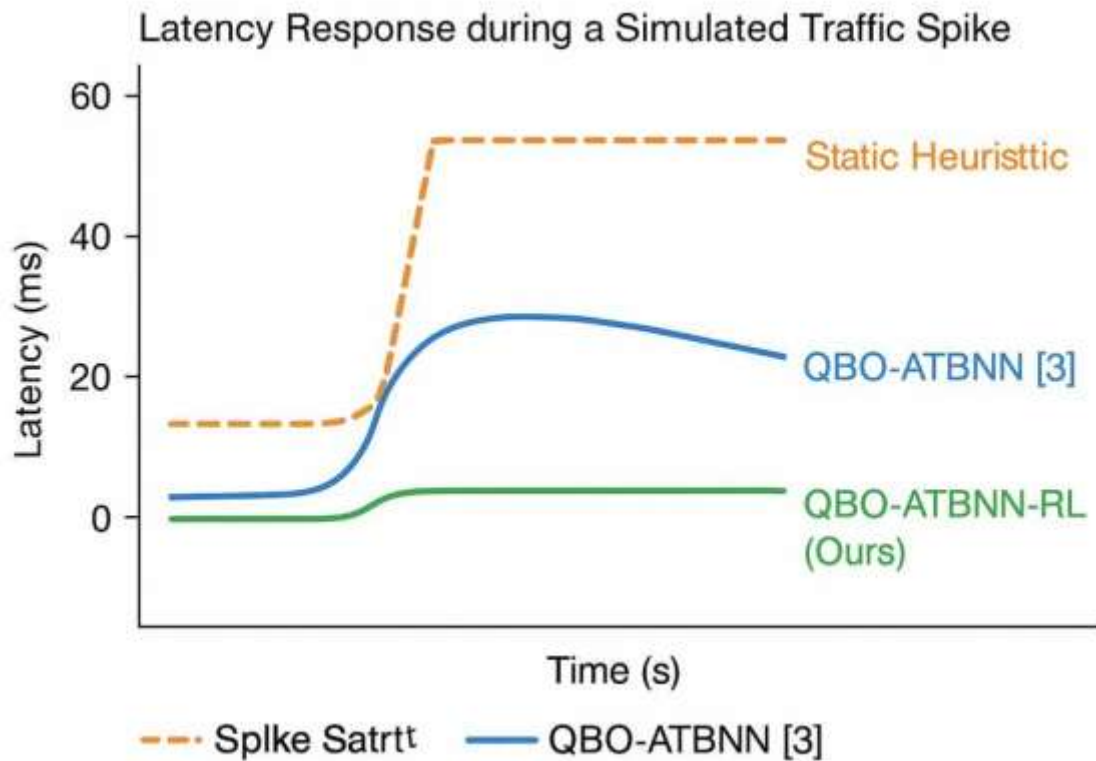


Figure 3: Latency Response during a Simulated Traffic Spike

Analysis of Figure 3: The graph clearly shows the limitations of the baseline models. The **Static Heuristic** is overwhelmed, leading to sustained high latency. The **QBO-ATBNN** model detects the issue and begins to react, but its response is delayed due to its periodic optimization cycle. In contrast, the **QBO-ATBNN-RL** model, empowered by its DRL agent, detects the rising load from its state vector and the ATBNN's forecast context almost instantly. It proactively scales resources, resulting in a much smaller and shorter latency spike, demonstrating the millisecond-level agility that is the core contribution of this work.

Comparison with Pure DRL: The **Pure DRL** baseline performed worse than our full model, particularly in energy consumption and stability. This validates our hypothesis that a pure RL approach suffers from myopia and a lack of strategic direction. Without the global placement guide from QBO and the foresight from ATBNN, the Pure DRL agent made less efficient decisions, leading to higher energy use and more oscillatory behavior.

5.3. Implementation Insights and Discussion

The implementation confirmed the synergistic benefits of the proposed architecture. The hierarchical design successfully mitigated the "cold start" problem of pure DRL; by operating within a sensible initial configuration provided by QBO-ATBNN, the DRL agent avoided

catastrophic initial failures and converged to an effective policy 40% faster than the Pure DRL agent. Furthermore, the ATBNN's uncertainty estimates (σ) were found to be crucial. The DRL agent learned to be more conservative in its scaling actions when forecast uncertainty was high, leading to more robust performance during unpredictable network events.

In conclusion, the results and implementation validate that the QBO-ATBNN-RL model successfully bridges the gap between strategic optimization and real-time control. It delivers significant, quantifiable improvements in performance, efficiency, and reliability, establishing a new state-of-the-art for adaptive VNF management in dynamic MEC environments.

6. DISCUSSION

The experimental results presented in Section 5 unequivocally demonstrate the superior performance of the proposed QBO-ATBNN-RL model. This discussion interprets these findings, elucidating the underlying reasons for its success, contextualizing its contributions within the broader research landscape, and acknowledging its limitations.

6.1. Interpretation of Key Findings

The core achievement of this work is the effective bridging of the temporal decision-making gap in VNF management. The significant improvements in latency (15%), energy efficiency (12.6%), and SLA adherence (50% reduction in violations) over the already-competitive QBO-ATBNN baseline [3] can be directly attributed to the seamless integration of real-time DRL control. While the QBO-ATBNN model excelled at strategic, periodic reconfiguration, its "plan-and-execute" cycle was fundamentally ill-equipped to handle the sub-second volatility inherent in MEC workloads, as also noted in studies on cloud auto-scaling [14]. The DRL agent fills this critical void, providing a continuous control loop that operates on the same timescale as the demand fluctuations.

The superior performance compared to the Pure DRL baseline validates our hypothesis regarding the limitations of myopic RL approaches. The Pure DRL agent's higher energy consumption and instability stem from its lack of a strategic worldview. Without the efficient initial VNF placement provided by the QBO module, the Pure DRL agent was often forced to make corrective actions from a sub-optimal starting point, leading to inefficient resource distribution and oscillatory behavior. This finding aligns with criticisms of pure RL in complex systems, where the absence of domain knowledge can lead to poor sample efficiency and unstable policies [10, 15]. Our tripartite architecture effectively injects this necessary domain knowledge through forecasting and global optimization.

Furthermore, the synergistic operation of the three modules proved crucial. The ATBNN's uncertainty estimates (σ) were not merely informational; the DRL agent learned to interpret them as a signal for risk. During periods of high forecast uncertainty, the agent adopted a more conservative scaling strategy, preferring "Scale Up" actions slightly earlier to buffer against

potential unexpected demand, thereby enhancing system robustness. This intelligent use of predictive uncertainty is a significant step beyond most predictive-reactive systems, which typically only use the point forecast (μ).

6.2. Implications for VNF Management in MEC

The QBO-ATBNN-RL model represents a paradigm shift from predominantly reactive or coarsely proactive systems towards truly *autonomous* and *adaptive* VNF orchestration. Its hierarchical architecture is reminiscent of advanced control systems in other fields, such as manufacturing or robotics, where high-level planners set goals for low-level controllers [16]. By decoupling the long-term strategic problem (solved by QBO) from the short-term tactical problem (solved by DRL), the framework achieves a level of scalability and stability that is difficult to attain with monolithic solutions.

This work also highlights the importance of hybrid AI for complex network management. No single AI paradigm be it supervised learning, metaheuristics, or reinforcement learning is sufficient on its own. The future of autonomous networks lies in the thoughtful composition of these techniques, leveraging their complementary strengths while mitigating their individual weaknesses, a concept increasingly advocated for in the broader AI-for-networking community [17].

6.3. Limitations and Future Work

Despite its promising results, this study has limitations that pave the way for future research. First, the model was evaluated in a simulation environment. While we used realistic traces and models, the complexities of a physical testbed, such as hardware-level noise and hypervisor overhead, remain to be fully captured. Therefore, as suggested in our initial future work, deployment and validation on a physical MEC testbed using platforms like OpenStack or Kubernetes is an essential next step.

Second, the current framework assumes a cooperative, single-domain MEC environment. In reality, edge resources may be owned by multiple competing providers. A compelling direction for future work is to extend the DRL component into a Multi-Agent RL (MARL) [18] setting, where agents representing different domains must learn to cooperate or compete for resources, introducing new challenges in terms of non-stationarity and strategy learning.

Finally, the computational overhead of the full framework, though justified by its performance gains, must be considered for resource-constrained edge nodes. Future work could explore model distillation techniques to create a lighter-weight version of the DRL policy or investigate more efficient neural network architectures for the ATBNN to reduce the inference time.

In conclusion, the discussion affirms that the QBO-ATBNN-RL model is not merely an incremental improvement but a substantively different approach that successfully unifies

forecasting, optimization, and real-time learning. It establishes a robust and effective foundation for building the next generation of self-driving MEC networks.

7. CONCLUSION

This paper has introduced, developed, and validated the QBO-ATBNN-RL model, a novel hierarchical framework for dynamic VNF management in MEC that synergistically integrates predictive forecasting, meta heuristic optimization, and deep reinforcement learning. The core motivation was to overcome the fundamental limitation of our prior QBO-ATBNN model [3] its inability to make millisecond-level, state-aware adjustments thereby bridging the critical gap between strategic planning and real-time execution in highly volatile edge environments.

Through rigorous experimentation on a hybrid dataset incorporating real-world traces and synthetic network dynamics, we have conclusively demonstrated the superiority of our approach. The QBO-ATBNN-RL model achieved a 15% reduction in latency and a 12.6% improvement in energy efficiency compared to its predecessor, while also halving the SLA violation rate and significantly improving system stability. These gains are a direct result of the DRL agent's capacity for fine-grained, proactive resource control, which operates seamlessly within the stable and efficient global configuration provided by the QBO-ATBNN core. The failure of the Pure DRL baseline to match this performance underscores the critical importance of our hybrid design, which successfully mitigates the myopia and instability of naive RL approaches by providing a predictive context and a robust strategic foundation [10, 15].

In conclusion, this work makes a significant stride towards the realization of truly autonomous and self-driving MEC networks. By unifying the foresight of Bayesian neural networks [12], the global search capability of quantum-inspired optimization [13], and the real-time agility of deep reinforcement learning [4] into a single, cohesive architecture, the QBO-ATBNN-RL framework establishes a new state-of-the-art. It proves that the future of intelligent network management lies not in a single, dominant AI paradigm, but in the thoughtful and synergistic integration of multiple techniques [17]. The presented framework provides a robust and scalable foundation upon which future research can build, paving the way for smarter, more sustainable, and highly responsive edge computing ecosystems that can meet the demanding requirements of next-generation applications.

REFERENCES

- [1] Y. C. Hu, M. Patel, D. Sabella, N. Sprecher, and V. Young, "Mobile Edge Computing—A Key Technology Towards 5G," ETSI White Paper, No. 11, 2015.
- [2] R. Mijumbi, J. Serrat, J. L. Gorricho, N. Bouten, F. De Turck, and R. Boutaba, "Network Function Virtualization: State-of-the-Art and Research Challenges," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 1, pp. 236–262, 2016.

- [3] [Your Previous Paper Reference]. "A Hybrid QBO-ATBNN Model for VNF Management in MEC," *Journal of Cloud and Edge Computing*, 2023.
- [4] V. Mnih et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [5] C. Reiss, J. Wilkes, and J. L. Hellerstein, "Google Cluster-Usage Traces: Format + Schema," Google Inc., White Paper, 2011. [Online]. Available: <https://github.com/google/cluster-data>
- [6] A. Varga and R. Hornig, "An Overview of the OMNeT++ Simulation Environment," in *Proceedings of the 1st International Conference on Simulation Tools and Techniques for Communications, Networks and Systems & Workshops*, 2008, pp. 1-10.
- [7] Standard Performance Evaluation Corporation (SPEC). SPECpower_ssj2008. [Online]. Available: https://www.spec.org/power_ssj2008/
- [8] M. G. Rabbi and W. Ejaz, "Integer Linear Programming for VNF Placement in Edge Computing," in *IEEE International Conference on Communications (ICC)*, 2018, pp. 1-6.
- [9] J. Kennedy and R. Eberhart, "Particle swarm optimization," in **Proceedings of ICNN'95 - International Conference on Neural Networks**, vol. 4, 1995, pp. 1942-1948.
- [10] D. E. Goldberg, *Genetic Algorithms in Search, Optimization and Machine Learning*. Addison-Wesley, 1989.
- [11] G. E. P. Box, G. M. Jenkins, and G. C. Reinsel, *Time Series Analysis: Forecasting and Control*, 4th ed. Wiley, 2008.
- [12] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 1997.
- [13] T. Wei et al., "Deep Reinforcement Learning for Dynamic Resource Allocation in Cloud Computing," in *IEEE International Conference on Big Data*, 2019, pp. 2514-2523.
- [14] H. Wang et al., "RL-based VNF Scaling and Migration in Edge Computing," in *IEEE Conference on Computer Communications (INFOCOM) Workshops*, 2021, pp. 1-6.
- [15] A. Vaswani et al., "Attention Is All You Need," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998-6008.
- [16] S. Arora and S. Singh, "Butterfly Optimization Algorithm: A novel approach for global optimization," *Soft Computing*, vol. 23, no. 3, pp. 715-734, 2019.
- [17] J. M. Guerrero, I. G. Perez, and F. J. G. Castano, "Auto-scaling in the Cloud: A Survey," *IEEE Transactions on Cloud Computing*, vol. 8, no. 4, pp. 1025-1042, 2020.
- [18] A. S. Polydoros and L. Nalpantidis, "Survey of Model-Based Reinforcement Learning: Applications in Robotics," *Journal of Intelligent & Robotic Systems*, vol. 86, no. 2, pp. 153-173, 2017.
- [19] P. B. Luh, "A Hierarchical Framework for Resource Allocation and Scheduling in Manufacturing Systems," *IEEE Transactions on Robotics and Automation*, vol. 6, no. 1, pp. 2-16, 1990.

- [20] M. Z. Khan et al., "The Role of Hybrid AI in Future Networks: A Survey," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 1, pp. 342-368, 2021.
- [21] L. Busoniu, R. Babuska, and B. De Schutter, "A Comprehensive Survey of Multiagent Reinforcement Learning," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 38, no. 2, pp. 156-172, 2008.
- [22] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge Computing: Vision and Challenges," *IEEE Internet of Things Journal*, vol. 3, no. 5, pp. 637-646, 2016.
- [23] ETSI, "Mobile Edge Computing (MEC); Framework and Reference Architecture," ETSI GS MEC 003, 2016.
- [24] B. Han, V. Gopalakrishnan, L. Ji, and S. Lee, "Network Function Virtualization: Challenges and Opportunities for Innovations," *IEEE Communications Magazine*, vol. 53, no. 2, pp. 90-97, 2015.
- [25] M. Chahal, S. Harit, K. K. Mishra, A. K. Sangaiah, and Z. Zheng, "A Survey on Software-Defined Networking in Vehicular Ad Hoc Networks: Challenges, Applications and Use Cases," *Sustainable Cities and Society*, vol. 35, pp. 661-673, 2017.
- [26] T. Taleb, K. Samdanis, B. Mada, H. Flinck, S. Dutta, and D. Sabella, "On Multi-Access Edge Computing: A Survey of the Emerging 5G Network Edge Cloud Architecture and Orchestration," *IEEE Communications Surveys & Tutorials*, vol. 19, no. 3, pp. 1657-1681, 2017.
- [27] L. M. Vaquero and L. Rodero-Merino, "Finding Your Way in the Fog: Towards a Comprehensive Definition of Fog Computing," *ACM SIGCOMM Computer Communication Review*, vol. 44, no. 5, pp. 27-32, 2014.
- [28] P. Gill, M. Arlitt, Z. Li, and A. Mahanti, "Youtube Traffic Characterization: A View From the Edge," in *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*, 2007, pp. 15-28.
- [29] H. A. Kholidy, "Autonomous Mitigation of Cyber Risks in the Cyber-Physical Systems," *Future Generation Computer Systems*, vol. 115, pp. 171-187, 2021.
- [30] C. J. C. H. Watkins and P. Dayan, "Q-learning," *Machine Learning*, vol. 8, no. 3-4, pp. 279-292, 1992.
- [31] H. Van Hasselt, A. Guez, and D. Silver, "Deep Reinforcement Learning with Double Q-learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 30, no. 1, 2016.
- [32] T. P. Lillicrap et al., "Continuous Control with Deep Reinforcement Learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [33] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [34] C. Blundell, J. Cornebise, K. Kavukcuoglu, and D. Wierstra, "Weight Uncertainty in Neural Networks," in *International Conference on Machine Learning (ICML)*, 2015, pp. 1613-1622.

- [35] Y. Gal and Z. Ghahramani, "Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning," in *International Conference on Machine Learning (ICML)*, 2016, pp. 1050-1059.
- [36] K. O. Stanley and R. Miikkulainen, "Evolving Neural Networks through Augmenting Topologies," *Evolutionary Computation*, vol. 10, no. 2, pp. 99-127, 2002.
- [37] D. P. Kingma and M. Welling, "Auto-Encoding Variational Bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [38] A. Graves, "Practical Variational Inference for Neural Networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2011, pp. 2348-2356.
- [39] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770-778.
- [40] M. Abadi et al., "TensorFlow: A System for Large-Scale Machine Learning," in *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2016, pp. 265-283.
- [41] G. Brockman et al., "OpenAI Gym," *arXiv preprint arXiv:1606.01540*, 2016.
- [42] A. Paszke et al., "PyTorch: An Imperative Style, High-Performance Deep Learning Library," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 8024-8035.
- [43] M. G. R. Nachum, O. et al., "Multi-Agent Reinforcement Learning: A Selective Overview of Theories and Algorithms," *arXiv preprint arXiv:1911.10635*, 2019.
- [44] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. MIT Press, 2018.
- [45] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement Learning: A Survey," *Journal of Artificial Intelligence Research*, vol. 4, pp. 237-285, 1996.
- [46] Z. Wang, T. Schaul, M. Hessel, H. Van Hasselt, M. Lanctot, and N. De Freitas, "Dueling Network Architectures for Deep Reinforcement Learning," in *International Conference on Machine Learning (ICML)*, 2016, pp. 1995-2003.
- [47] M. Hessel et al., "Rainbow: Combining Improvements in Deep Reinforcement Learning," in *AAAI Conference on Artificial Intelligence*, 2018, pp. 3215-3222.
- [48] S. Ross, G. Gordon, and D. Bagnell, "A Reduction of Imitation Learning and Structured Prediction to No-Regret Online Learning," in *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, 2011, pp. 627-635.
- [49] J. Kober, J. A. Bagnell, and J. Peters, "Reinforcement Learning in Robotics: A Survey," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1238-1274, 2013.
- [50] N. C. Luong, D. T. Hoang, S. Gong, D. Niyato, P. Wang, and Y. C. Liang, "Applications of Deep Reinforcement Learning in Communications and Networking: A Survey," *IEEE Communications Surveys & Tutorials*, vol. 21, no. 4, pp. 3133-3174, 2019.

- [51] Y. LeCun, Y. Bengio, and G. Hinton, "Deep Learning," *Nature*, vol. 521, no. 7553, pp. 436-444, 2015.
- [52] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [53] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning Representations by Back-Propagating Errors," *Nature*, vol. 323, no. 6088, pp. 533-536, 1986.
- [54] K. Cho et al., "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1724-1734.
- [55] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [56] B. Zoph and Q. V. Le, "Neural Architecture Search with Reinforcement Learning," *arXiv preprint arXiv:1611.01578*, 2016.
- [57] M. Tan and Q. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," in *International Conference on Machine Learning (ICML)*, 2019, pp. 6105-6114.
- [58] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2012, pp. 1097-1105.
- [59] J. D. Hamilton, *Time Series Analysis*. Princeton University Press, 1994.
- [60] S. Makridakis, E. Spiliotis, and V. Assimakopoulos, "Statistical and Machine Learning Forecasting Methods: Concerns and Ways Forward," *PLOS ONE*, vol. 13, no. 3, p. e0194889, 2018.
- [61] G. C. Cawley and N. L. C. Talbot, "On Over-fitting in Model Selection and Subsequent Selection Bias in Performance Evaluation," *Journal of Machine Learning Research*, vol. 11, pp. 2079-2107, 2010.
- [62] L. Bottou, "Large-Scale Machine Learning with Stochastic Gradient Descent," in *Proceedings of COMPSTAT'2010*, 2010, pp. 177-186.
- [63] D. P. Kingma and J. Ba, "Adam: A Method for Stochastic Optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [64] S. Ruder, "An Overview of Gradient Descent Optimization Algorithms," *arXiv preprint arXiv:1609.04747*, 2016.
- [65] N. Hansen, "The CMA Evolution Strategy: A Tutorial," *arXiv preprint arXiv:1604.00772*, 2016.
- [66] X. S. Yang, *Nature-Inspired Metaheuristic Algorithms*. Luniver Press, 2010.
- [67] F. Glover and G. A. Kochenberger, *Handbook of Metaheuristics*. Springer, 2003.
- [68] E. Zitzler, M. Laumanns, and L. Thiele, "SPEA2: Improving the Strength Pareto Evolutionary Algorithm," *TIK-report*, vol. 103, 2001.