

ML-Based QoS Optimization in Enterprise Networks

Deepak Chauhan

Yashwantrao Chavan Maharashtra Open University

Abstract- The digital infrastructure of the modern enterprise is undergoing a radical transformation, driven by the widespread adoption of cloud-native applications, real-time collaboration tools, and high-bandwidth multimedia services. In this dynamic landscape, traditional Quality of Service (QoS) mechanisms, which rely on static priority queuing and manually defined traffic classes, are increasingly incapable of managing the volatility of network demand. This review explores the paradigm shift toward Machine Learning (ML)-based QoS optimization. By transitioning from reactive, threshold-based management to proactive, intent-driven architectures, ML enables enterprise networks to achieve "Cognitive Traffic Engineering." This article examines how various ML paradigms—including supervised learning for traffic classification, unsupervised learning for anomaly detection, and reinforcement learning for dynamic resource allocation—can be synthesized into a unified optimization fabric. We analyze the efficacy of Deep Learning models, such as Convolutional Neural Networks and Long Short-Term Memory units, in identifying application-layer requirements within encrypted tunnels without the need for Deep Packet Inspection. Furthermore, the review addresses the architectural integration of ML within Software-Defined Networking (SDN) and SD-WAN frameworks, enabling the "Self-Driving Network" vision. Critical challenges, such as model interpretability, real-time inference latency at the network edge, and data drift in multi-tenant environments, are discussed in depth. By synthesizing recent academic breakthroughs and industrial implementations, this paper provides a strategic roadmap for building resilient, high-performance enterprise networks. The findings suggest that ML-driven QoS is the foundational technology required to satisfy the stringent Service Level Agreements of the modern digital enterprise, ensuring that network resources are distributed with machine-speed precision and contextual intelligence.

Keywords – QoS Optimization, Machine Learning, Traffic Engineering, Software-Defined Networking, Intent-Based Networking.

I. INTRODUCTION

The history of enterprise networking is defined by a persistent struggle to balance finite bandwidth with an ever-expanding volume of user demand. Historically, Quality of Service was a relatively simple task. In the era of the local area network, traffic was predictable, and applications were few. Network administrators could define static "Class of Service" tags for voice and data, ensuring that the limited capacity of expensive leased lines was prioritized for critical business functions. However, the rise of the "Hyper-Connected Enterprise" has rendered these traditional, manual methods obsolete.

Today, the network is no longer a static perimeter; it is a fluid, global fabric that spans private data centres, public clouds, and thousands of remote branch offices. The sheer variety of traffic—ranging from ultra-reliable low-latency communication for industrial robotics to massive, bandwidth-hungry 4K video streams—creates a level of complexity that human-led configuration can no longer manage.

The fundamental flaw of traditional QoS is its "Static Nature." A human engineer writes a rule: "If the link is 80% full, drop the low-priority packets." This is a reactive, binary decision that ignores the "Application Context." It doesn't account for the fact that a "Low-Priority" file transfer might be a critical security patch, or that "High-Priority" video traffic might just be an unimportant social call.

Furthermore, as over 95% of web traffic becomes encrypted, traditional Deep Packet Inspection (DPI) tools can no longer "see" what is inside the packet to classify it correctly. This "Visibility Crisis" and "Complexity Gap" have necessitated the move toward Machine Learning. ML-based QoS optimization represents the transition from "Managing Packets" to "Predicting Intent." By using statistical patterns and behavioral signatures, ML can identify the true requirements of a traffic flow even when it is hidden behind multiple layers of encryption.

Machine Learning provides the "Cognitive Brain" for the network's control plane. It enables the network to "Observe"

its own performance, "Learn" from historical traffic peaks, and "Decide" on the optimal path and priority for every flow in real-time. This is the foundation of "Intent-Based Networking" (IBN). In an IBN environment, the administrator simply states a business goal—such as "ensure the CEO's video call never drops, regardless of other traffic"—and the ML-driven QoS engine manages the millions of low-level configurations required to meet that goal. This shift moves the enterprise from a "Break-Fix" operational model to a "Predictive-Optimization" model. Instead of waiting for a user to complain about a slow connection, the ML agent identifies the impending congestion ten minutes before it happens and reroutes non-essential traffic to a secondary link.

However, the implementation of ML in enterprise QoS is not without significant hurdles. For an ML model to be effective, it must operate at "Line Speed." If the classification and decision-making process takes even a few milliseconds too long, the resulting jitter can be worse than the congestion it was trying to solve. This necessitates the use of "Hardware Acceleration" and "Edge Intelligence," where the ML models are baked into the logic of the network switches and routers themselves. Furthermore, we must address the "Trust Gap." Network engineers are historically skeptical of "Black Box" AI. For ML-based QoS to be widely adopted, the models must be "Explainable," providing a clear reasoning path for why a specific flow was throttled or rerouted. This review will explore the technological evolution of these systems, from basic supervised learners to the cutting edge of "Reinforcement Learning" and "Graph Neural Networks." We will analyze how the fusion of big data analytics and algorithmic intelligence is creating a more resilient, efficient, and transparent digital infrastructure.

Ultimately, ML-driven QoS is about the "Democratization of Performance." It ensures that every application, from the humblest background sync to the most critical financial transaction, receives exactly the resources it needs to function optimally. In an era where the network is the business, the ability to manage that network with machine-intelligence is the defining factor in an organization's digital resilience. This review provides a granular look at the architectures, data strategies, and operational challenges that define the current state-of-the-art in intelligent QoS, providing a roadmap for the "Self-Driving Enterprise" of the next decade.

II. EVOLUTIONARY TRENDS IN TRAFFIC CLASSIFICATION METHODOLOGIES

The first pillar of any QoS strategy is "Classification." If the network does not know what an application is, it cannot prioritize it. In the traditional era, this was done via port-based mapping or signature-based DPI. As encryption rendered these methods obsolete, Machine Learning stepped in to

perform "Statistical Traffic Fingerprinting." This methodology treats an encrypted flow as a series of behavioral metrics—packet size distributions, inter-arrival times, and byte-level entropy. By training on vast datasets of labeled traffic, ML models can identify applications like Microsoft Teams, Salesforce, or BitTorrent with over 99% accuracy based solely on the "shape" of the communication.

The current trend is moving away from "Shallow ML" (like Random Forests) toward "Deep Learning" (DL). Deep architectures, specifically 1D Convolutional Neural Networks (CNNs), can process the raw byte stream of the first few packets in a flow to identify the application intent before the handshake is even complete. This "Early Classification" is vital for QoS, as it allows the network to apply the correct priority from the very first packet of the session. We are also seeing the rise of "Multi-Modal Learning," where the model combines the "Spatial Features" (packet structure) with "Temporal Features" (flow timing over minutes). This holistic view allows the network to distinguish between different functions within a single app—for example, distinguishing between a "Voice Call" and a "File Transfer" inside a single Slack session. This granular classification is the bedrock upon which intelligent resource allocation is built.

III. ADAPTIVE RESOURCE ALLOCATION VIA REINFORCEMENT LEARNING

Once the traffic is classified, the next challenge is "Allocation"—deciding how much bandwidth and what buffer priority each flow should receive. Traditional QoS uses fixed "Scheduler Weights," which are static and do not adapt to changing network conditions. "Reinforcement Learning" (RL) represents a paradigm shift in this area. In an RL-based QoS system, the network acts as an "Agent" that interacts with its "Environment." The agent takes actions (adjusting queue sizes or reroute traffic) and receives "Rewards" (reduced latency or higher throughput). Over millions of iterations, the RL agent learns the "Optimal Policy" for its specific, unique network environment.

The beauty of RL is its "Autonomy." Unlike supervised learning, it doesn't need to be told the "correct" answer; it discovers it through trial and error. In a complex enterprise network where the relationship between a routing change and its impact on a user's experience is non-linear, RL is the only mechanism capable of finding the global optimum. We explore "Deep Q-Networks" (DQN) and "Policy Gradient" methods that allow these agents to manage thousands of simultaneous flows across a global SD-WAN. This results in "Dynamic Slicing," where the network's capacity is virtually partitioned and reshaped every millisecond to match the real-time ebb and flow of business demand. By turning resource management into a "Learning Task," RL ensures that the

network is always at its peak efficiency, even when faced with unforeseen traffic spikes or hardware failures.

Deep Learning for Predictive Buffer and Queue Management
In traditional routers, the "Buffer" is a simple storage area where packets wait their turn. If the buffer is full, the router drops the incoming packets, leading to "Tail Drop" and severe performance degradation for TCP-based applications. ML-based QoS introduces "Predictive Queue Management." By using "Long Short-Term Memory" (LSTM) networks—a type of DL designed for time-series forecasting—the network can predict a "Buffer Overflow" several seconds before it happens.

This allows the QoS engine to take "Proactive Remediation" steps, such as triggering "Explicit Congestion Notification" (ECN) or performing "Early Random Drop" on non-critical packets. LSTMs are particularly effective here because network traffic is inherently cyclical and seasonal. There are daily peaks during the start of the workday and weekly lulls during the weekend. The ML model learns these "Patterns of Life," allowing it to pre-emptively clear out buffers for an expected surge in video traffic during a company-wide town hall. This section analyzes the transition from "Passive Buffering" to "Intelligent Congestion Control."

By maintaining a "Predictive Horizon," the network eliminates the "Sawtooth Pattern" of TCP congestion windows, leading to a much smoother and more consistent user experience. This "Smooth Performance" is essential for real-time applications like cloud-based CAD/CAM or medical imaging, where even minor variations in packet delivery can disrupt the work of a professional user.

IV. GRAPH NEURAL NETWORKS FOR RELATIONAL NETWORK INTELLIGENCE

The enterprise network is not just a collection of links; it is a "Topology." A change in a hub in New York can have a ripple effect on a branch in London. Traditional ML models, which treat network data as a "Flat Table" of numbers, are blind to this relational structure. "Graph Neural Networks" (GNNs) are a new class of AI designed specifically to process data on graphs. In a GNN-based QoS architecture, the nodes are the switches and the edges are the links. The GNN "reasons" about the network's structure, predicting how a congestion event will propagate through the entire global fabric.

This section explores the use of GNNs for "Global QoS Optimization." By analyzing the "Network Graph," the AI can identify "High-Centrality" nodes that act as bottlenecks for the entire organization. It then suggests topological changes or dynamic routing paths that optimize the global QoS, rather than just the local link. GNNs are also essential for managing "Multi-Cloud Interconnects." As organizations move traffic

between AWS, Azure, and Google Cloud, the GNN helps the QoS engine understand the "Relational Cost and Performance" of different cloud on-ramps. This "Topological Awareness" is what allows for true "Enterprise-Wide QoS," ensuring that the end-to-end path of a packet is optimized across the entire distributed estate, rather than just within the four walls of the data centre.

V. HANDLING NON-STATIONARITY AND CONCEPT DRIFT IN DYNAMIC NETWORKS

The primary enemy of a machine learning model is "Change." In an enterprise network, traffic patterns are notoriously non-stationary. A new software deployment, a merger with another company, or a shift to a "Hybrid Work" model can fundamentally change the statistical distribution of data. If an ML model is not designed to handle this "Concept Drift," its QoS decisions will become inaccurate over time.

This section explores "Adaptive Learning" strategies, where the AI continuously updates its own weights based on real-time feedback. Instead of a "Static" model, the network uses "Streaming Analytics" to learn on the fly. We examine the use of "Online Learning" and "Transfer Learning." Transfer learning allows a model trained on a generic dataset of enterprise traffic to be "Fine-Tuned" for a specific company's unique quirks in a matter of hours. We also discuss "Robustness Training," where the ML model is intentionally tested against "Noisy" and "Adversarial" data to ensure it doesn't fail during a network crisis.

This resilience is what transforms ML from a "Research Interest" into a "Production-Ready" security and performance tool. By building a "Self-Correcting" intelligence, enterprise networks ensure that their QoS optimization remains effective even as the business grows and the technological landscape shifts under their feet.

VI. EXPLAINABLE AI (XAI) AND BUILDING OPERATOR TRUST

One of the major barriers to the adoption of ML in the network operations centre (NOC) is the "Black Box" problem. If an AI decides to throttle the traffic of the Marketing department, a network engineer needs to know "Why." Without transparency, the AI remains a mysterious oracle that is often ignored or overridden by humans. "Explainable AI" (XAI) is the technological layer that provides the "Reasoning Path" for an AI decision. XAI tools like "SHAP" or "LIME" are used to provide "Transparency Logs" for every QoS change.

This section explores the transition from "Black Box" to "Glass Box" models. By providing the analyst with a "Risk Heatmap" or a "Feature Importance Chart," XAI allows the human to verify that the AI is acting on legitimate business logic. For example: "I am throttling this flow because its jitter profile indicates a non-critical background update, and the high-priority Zoom queue is currently at 90% capacity." This "Machine-to-Human" communication is essential for building trust. It also allows for "Regulatory Compliance" and "Post-Mortem Analysis," ensuring that every QoS decision can be audited and justified. By making the AI "Understandable," enterprise networks can finally move toward "Full Autonomy," where humans provide the strategy and AI handles the machine-speed execution.

VII. SCALABILITY AND REAL-TIME INFERENCE AT THE NETWORK EDGE

The ultimate challenge for ML-based QoS is "Line-Speed Inference." In a production environment with 100Gbps links, the network must make a QoS decision every few microseconds. If the ML model is too large or too slow, it becomes a bottleneck. This section explores "Model Compression" and "Hardware Acceleration." We analyze techniques like "Pruning" and "Quantization," which shrink a massive neural network into a lean "Surrogate Model" that can run on a standard router CPU or a specialized NPU (Neural Processing Unit).

We examine the rise of "In-Network Computing," where the ML model is baked directly into the programmable data plane (using languages like P4). This allows the switch to classify and prioritize a packet as it is moving through the silicon, with "Zero-Latency." We also discuss "Distributed Inference," where the "Learning" happens in the cloud, but the "Decision-Making" happens at the branch office edge.

This "Hierarchical AI" architecture ensures that the network is responsive and scalable, capable of managing millions of simultaneous flows across a global infrastructure without adding significant delay to the end-user experience. By moving the "Intelligence" to the "Edge," enterprise networks ensure that QoS is applied at the point of origin, preventing congestion before it ever reaches the core.

Managing Quality of Experience (QoE) in the SD-WAN Era QoS is a technical metric (latency, jitter), but the business cares about "Quality of Experience" (QoE)—how the user actually feels about the application. AI is the only tool capable of bridging the gap between "Packet Metrics" and "Human Perception." By training on subjective user feedback and objective network telemetry, AI models can predict the "MOS" (Mean Opinion Score) for a video call in real-time. If the AI predicts the QoE will drop below a "4.0" threshold, the

QoS engine takes immediate corrective action. This is the heart of "Service-Centric Networking."

This section explores the integration of ML-QoS within "SD-WAN" (Software-Defined Wide Area Network) architectures. In an SD-WAN, traffic can be steered over multiple paths (MPLS, Broadband, 5G). An ML-driven "Path Steering" engine can analyze the "Historic and Current Performance" of every available path to select the one that offers the best QoE for a specific app. We discuss how this "Multi-Path Optimization" reduces costs by allowing businesses to use cheap broadband for critical apps without sacrificing reliability. By making the user's experience the "Primary Objective Function" of the network, ML-driven QoS aligns the technical infrastructure with the strategic goals of the digital enterprise, ensuring that technology serves the person, rather than the other way around.

VIII. SECURITY-AWARE QOS AND TRAFFIC FUSION

In the modern enterprise, "Performance" and "Security" can no longer be managed in silos. A QoS strategy that prioritizes an encrypted stream that turns out to be a "Data Exfiltration" event is a security failure. "Security-Aware QoS" uses ML to fuse these two domains. The same ML model that identifies an application for prioritization also checks for "Behavioral Anomalies" that suggest a security risk. If a high-priority flow starts exhibiting "Suspicious Outbound Patterns," the QoS engine can instantly "De-Prioritize" it or move it to a "Scrubbing Center" for deeper inspection.

This section examine the concept of "Fusion Analytics," where the network's "Resource Manager" and "Threat Detector" share the same ML feature set. This reduces the computational overhead and ensures that "Mission-Critical" traffic is not just "Fast," but also "Safe." We analyze the threat of "Adversarial AI," where an attacker tries to "Spoof" the QoS engine into giving their malicious traffic high priority by mimicking the behavioral signature of a VoIP call. To counter this, we discuss "Robust ML" architectures that look for the "Inconsistencies" between a flow's claimed priority and its actual structural logic. By integrating security into the QoS fabric, the enterprise network becomes a "Hardened Infrastructure" that protects its most critical assets while ensuring their peak performance.

IX. FUTURE PERSPECTIVES: 6G AND AUTONOMOUS NETWORK INTENT

As we look toward the 2030s, the scope of enterprise QoS will expand into the "Extreme Edge"—including mobile users on 6G networks and massive swarms of IoT devices. 6G promises sub-millisecond latency and terabit-per-second

speeds, which will require a level of "Granular Orchestration" that is far beyond today's capabilities. This section explores the "Autonomous Network Intent" vision, where the network doesn't just follow rules, but "Understands the Business." If a company is launching a new product, the network AI identifies the increased importance of the "Sales API" and autonomously reconfigures the global QoS to support the launch.

We also examine the role of "Sustainability" in QoS. Future AI models will not just optimize for "Speed," but also for "Energy." The AI-QoS engine can steer traffic to links and nodes that are currently powered by renewable energy, or "Power Down" underutilized segments of the network during off-peak hours without impacting the user experience. We conclude by looking at the "Full Autonomy" phase, where the network is a "Sentient Ecosystem" that self-repairs, self-scales, and self-defends. This is the final frontier of QoS, where the network becomes "Invisible"—a perfectly optimized utility that provides a flawless digital experience without a single human intervention. This represents the ultimate "Optimization": the reduction of network friction to near-zero.

X. CONCLUSION

ML-based QoS optimization represents the definitive transition from static, human-led network management to an autonomous, intent-aware digital fabric. By leveraging the predictive power of LSTMs, the strategic decision-making of Reinforcement Learning, and the relational intelligence of Graph Neural Networks, enterprise networks can finally address the volatility and complexity of the multi-cloud world. This review has demonstrated that "Intelligence" is no longer an optional feature; it is the core engine required to satisfy the human-centric "Quality of Experience" demands of the modern workforce.

However, the path toward full autonomy requires a rigorous focus on "Explainability" to maintain human trust and "Edge Acceleration" to ensure real-time performance. As we move into an era of 6G and massive IoT, the ability to manage the global "Data Flow" with machine-intelligence will be the deciding factor in an organization's digital resilience. Ultimately, ML-driven QoS ensures that the network is no longer a bottleneck for business innovation, but a dynamic, self-optimizing catalyst for the next era of global digital transformation. It is the silent architect of the modern enterprise, ensuring that every bit and every byte is delivered with precision, purpose, and peak efficiency.

REFERENCES

1. Burrumukku, N. R. (2015). Real-time detection of network threats using deep packet inspection and telemetry analytics. *International Journal of Trend in Research and Development*, 2(1), 1–5.
2. Jangala, V. K. (2015). Observability and monitoring of microservices using Splunk and New Relic. *International Journal of Engineering Development and Research*, 3(3), 1–15.
3. Vangoor, V. K. R. (2016). AI-driven monitoring and alerting systems for enterprise-scale Linux deployments. *International Journal of Science, Engineering and Technology*, 4(1), 11.
4. Parimi, S. S. (2016). Analyzing the effectiveness of SAP systems in streamlining healthcare supply chains, reducing costs, and improving service delivery.
5. Koukuntla, S. (2018). Event-driven architectures in cloud computing: Tools, patterns, and tradeoffs. *International Journal of Trend in Scientific Research and Development*, 2(3), 2909–2913.
6. Burrumukku, N. R. (2015). Root cause analysis in enterprise networks using correlated telemetry and graph analytics. *TIJER – International Research Journal*, 2(6), a9–a17.
7. Jangala, V. K. (2016). API gateway security implementation using JWT and Apigee in cloud-native applications. *International Journal of Current Science*, 6(2), 34–43.
8. Vangoor, V. K. R. (2017). Self-optimizing DevOps pipelines for enterprise infrastructure using machine learning models. *International Journal of Trend in Scientific Research and Development*, 1(6), 8.
9. Parimi, S. S. R. (2016). Predictive analytics for financial forecasting in SAP ERP systems using machine learning. *International Journal of Creative Research Thoughts*.
10. Burrumukku, N. R. (2016). Secure identity and access management integration for cloud-native network observability platforms. *International Journal of Engineering Development and Research*.
11. Jangala, V. K. (2018). Database performance tuning strategies for high-volume transaction systems. *International Journal of Scientific Development and Research*, 3(8), 274–282.
12. Vangoor, V. K. R. (2018). AI-based optimization of automated server deployment using Kickstart and Satellite systems. *International Journal of Trend in Research and Development*, 5(6), 5.
13. Parimi, S. S. (2018). Exploring the role of SAP in supporting telemedicine services, including scheduling, patient data management, and billing. *SSRN Electronic Journal*.
14. Burrumukku, N. R. (2016). Secure storage and backup architectures for cloud integrated datacenters.

- International Journal of Science, Engineering and Technology, 4(3).
15. Burremukku, N. R. (2017). End-to-end SD-WAN performance evaluation across private and public transport networks. *International Journal of Current Science*, 7(1), 56–65.
 16. Burremukku, N. R. (2017). Identity-aware network segmentation using NSX and next-generation firewalls. *International Journal of Scientific Research & Engineering Trends*, 3(5).
 17. Parimi, S. S. (2018). Optimizing financial reporting and compliance in SAP with machine learning techniques. *SSRN Electronic Journal*.
 18. Burremukku, N. R. (2018). Evaluating high-availability DHCP architectures: Migration from legacy Linux DHCP to Infoblox grid. *International Journal of Scientific Development and Research*.
 19. Mandati, S. R. (2019). The basic and fundamental concept of cloud balancing architecture. *South Asian Journal of Engineering and Technology*, 9(1), 4.