

QUANTUM QUEUEING THEORY: MODELING AND ANALYSIS OF QUANTUM COMPUTING-BASED SERVICE SYSTEMS: A RESEARCH ROADMAP

¹Dr.T.Vengatesh, ¹Kalpana Devarajan, ²Dr.K. Prabhavathi, ³Dr.Brinda Halambi, ⁴Saint Jesudoss.S, ⁵Dr.S.S.Ananthan, ⁶Dr.B.Anbuselvan, ⁷Dr.L.Jerlin Rubini

1*(Coressponding Author) Assistant Professor, Department of Computer Science, Government Arts and Science College, Veerapandi, Theni, Tamilnadu, India. Email ID: venkibiotinix@gmail.com

¹Associate Professor, Department of Mathematics, KIT- Kalaigarkarananidhi Institute of Technology, Coimbatore-641402, Tamilnadu, India. Email ID: devkalpu@gmail.com

²Assistant Professor(Selection grade), Department of Mathematics, Bannari Amman Institute of Technology, Sathyamangalam -638401, Tamilnadu,India. Mail ID: prabhavathik@bitsathy.ac.in

³Associate Professor, Department of Mathematics, School of Applied Sciences, Reva University, Bangalore, India. EMail ID: brindahalambi@gmail.com

⁴Assistant Professor,Department of CSE, Rajiv Gandhi College of Engineering and Technology, Puducherry, India , EMail ID: saint.2k5@gmail.com

⁵Associate Professor, Department of Mathematics, Erode Sengunthar Engineering College, Thudupathi, Perundurai(TK) 638 057, Erode, Tamil Nadu, India. [EMail ID: ananthmathsesecc@gmail.com](mailto:ananthmathsesecc@gmail.com)

⁶Assistant Professor, Department of Computer Science, Government Arts and Science College, Veerapandi, Theni, Tamilnadu, India. Email ID: selvananbume@gmail.com

⁷Assistant Professor, Department of Computer Science, Government Arts and Science College, Veerapandi, Theni, Tamilnadu, India. Email ID: jel.jerlin@gmail.com

**^{1*} Corresponding Author:^{1*} Dr.T.VENGATESH,
Email ID: venkibiotinix@gmail.com**

ABSTRACT:

The convergence of quantum computing and classical service system modeling presents a transformative opportunity. This paper proposes the formalization of Quantum Queueing Theory (QQT), a novel framework for modeling, analyzing, and optimizing service systems where quantum processors act as servers, quantum algorithms constitute service tasks, and quantum communication channels facilitate arrivals and departures. We outline the fundamental challenges posed by quantum mechanics (superposition, entanglement, measurement, decoherence) to classical queueing paradigms. Key research directions include defining quantum analogues of arrival processes, service disciplines, and performance metrics (quantum fidelity, task success probability, decoherence-limited waiting time). We explore modeling approaches leveraging quantum stochastic processes, quantum walks, and modified Lindblad master equations. The paper details critical areas for future research: hybrid quantum-classical queueing networks, resource allocation under decoherence, stability analysis in the quantum regime, and the development of quantum-aware scheduling policies. QQT is poised to become essential for designing efficient and scalable quantum computing data centers, quantum cloud services, and integrated quantum-classical computing infrastructures.

Keywords: Quantum Queueing Theory (QQT), Quantum Computing, Service Systems, Quantum Servers, Decoherence, Quantum Scheduling, Quantum Stochastic Processes, Quantum Cloud Services, Hybrid Quantum-Classical Systems.

1.INTRODUCTION:

The burgeoning field of quantum computing promises unprecedented computational power, necessitating the design of efficient supporting service infrastructures like data centers and cloud platforms. However, classical queueing theory, the bedrock for modeling and optimizing traditional service systems, faces fundamental limitations when applied to quantum environments due to the unique characteristics of quantum mechanics superposition, entanglement, measurement, and decoherence. This paper formally introduces Quantum Queueing Theory (QQT) as a novel framework specifically designed to model, analyze, and optimize service systems where quantum processors act as servers, quantum algorithms constitute service tasks, and quantum communication channels manage arrivals and departures. We outline the core challenges in translating classical queueing paradigms to the quantum realm and propose key research directions, including defining quantum arrival processes, service disciplines, and performance metrics (e.g., fidelity, success probability, decoherence-limited waiting time), leveraging tools like quantum stochastic processes and modified master equations. QQT provides the essential theoretical foundation for designing scalable and efficient quantum computing service infrastructures.

The emergence of Quantum Computing as a Service (QCaaS) and dedicated quantum data centers necessitates efficient management of hybrid quantum-classical workloads, highlighting a critical need for accurate performance prediction and resource allocation. However, classical queueing theory, foundational to classical data centers, faces significant limitations in this quantum context. Its core assumptions—deterministic service times, classical state representation, and memoryless processes—break down due to uniquely quantum phenomena such as superposition, entanglement, probabilistic measurement outcomes, and decoherence. Furthermore, standard performance metrics like throughput must be augmented with quantum-specific measures, particularly fidelity or success probability. This necessitates a new *quantum service paradigm*, where tasks (e.g., VQE, QPE, or segments of Shor's algorithm) arrive at quantum processing units (QPUs), demanding specific resources like qubits, gate operations, and coherence time, often involving quantum communication for task submission and result retrieval. To address these challenges, this paper proposes the framework of Quantum Queueing Theory (QQT). Our objective is to define QQT's scope, articulate its inherent challenges stemming from quantum mechanics, explore potential modeling approaches, and identify critical research questions, thereby motivating its necessity for the practical deployment and scaling of

quantum computing. The paper is structured to first establish this foundation before delving into specific QQT models, analysis, and future directions.

2.FUNDAMENTAL CONCEPTS & CHALLENGES

Quantum Queueing Theory (QQT) fundamentally reimagines service system modeling by placing uniquely quantum mechanical elements at its core. The central concept involves quantum servers – Quantum Processing Units (QPUs) – which execute quantum tasks (e.g., VQE, QPE, Shor's algorithm segments). These tasks arrive via quantum communication channels and demand specific quantum resources such as qubit count, gate operations, and crucially, coherence time. This quantum service paradigm inherently diverges from classical queueing models due to the irreducible characteristics of quantum mechanics: superposition allows tasks or server states to exist in multiple configurations simultaneously, entanglement creates complex, non-local correlations between tasks or server resources, quantum measurement collapses the state probabilistically upon task completion or observation, and decoherence imposes a finite, stochastic lifespan on quantum information. These phenomena collectively shatter the foundational assumptions of classical queueing theory. Deterministic service times become untenable due to the probabilistic nature of quantum operations and measurement outcomes. The requirement for a classical state representation fails as quantum systems inherently possess non-classical correlations and superposition. Memoryless processes are violated due to entanglement and the time-dependent decay caused by decoherence. Consequently, standard performance metrics like throughput and latency must be augmented with intrinsically quantum metrics, primarily fidelity (measuring the accuracy of the output state) and task success probability, alongside decoherence-limited waiting time reflecting the urgency imposed by qubit instability. This necessitates entirely new modeling approaches for defining quantum arrival processes (potentially involving entangled task streams), quantum service disciplines (scheduling entangled or superpositioned tasks), and capturing the complex interplay between task execution and the relentless decay of quantum information. Furthermore, the practical context of hybrid quantum-classical workloads introduces additional complexity, requiring QQT to seamlessly integrate classical and quantum queueing dynamics within a single framework. Addressing these challenges requires novel theoretical tools, potentially leveraging quantum stochastic processes, quantum walks, or modified Lindblad master equations, to accurately model and analyze the performance and stability of quantum computing service infrastructures.

3.MODELING FRAMEWORKS FOR QUANTUM QUEUEING THEORY

Quantum Queueing Theory (QQT) requires new mathematical frameworks that explicitly incorporate quantum mechanical phenomena while still adhering to queueing-theoretic principles. Three primary approaches serve as the foundational pillars for this field.

3.1. Quantum Stochastic Processes

Generalizing classical Poisson processes, quantum arrival streams may exhibit entanglement between tasks or superpositional arrivals. The system state is modeled as a density matrix ρ that evolves under quantum stochastic differential equations:

$$d\rho = -\hbar i[H, \rho]dt + k \sum (L_k \rho L_k^\dagger - 2\{L_k^\dagger L_k, \rho\})dt + J(\rho)dN_t$$

In this equation, H is the system Hamiltonian, the operators L_k model decoherence and dephasing, and $J(\rho)$ captures quantum jumps that result from task arrivals or departures with the counting process N_t .

3.2. Quantum Walks for Discrete-State Systems

For queue systems with discrete states, quantum walks on graphs model superpositional queue occupations. The state vector $|\psi(t)\rangle \in H_Q \otimes H_S$ evolves according to the following unitary operation:

$$|\psi(t+1)\rangle = U \cdot (S \otimes C) |\psi(t)\rangle$$

Here, H_Q is the queue occupancy space, H_S is the server state space, U is a unitary service operator, S is a shift operator for arrivals and departures, and C is a quantum coin operator that governs routing decisions.

3.3. Modified Lindblad Master Equations

To integrate decoherence during waiting times, extended **Lindblad equations** are used to track task-resource entanglement. The evolution is described by:

$$\rho' = -\hbar i[H_{\text{queue}}, \rho] + \sum_{i=1}^M \gamma_i (A_i \rho A_i^\dagger - 2\{A_i^\dagger A_i, \rho\})$$

In this equation, H_{queue} governs the coherent queue dynamics, while the jump operators A_i model various processes:

- **Task arrivals** with rate λ : $A_{\text{arr}} = \lambda \sum_n |n+1\rangle\langle n| \otimes \sigma_{\text{arr}}$ **Error! Filename not specified.**
- **Decoherence** during waiting: $A_{\text{dec}}(k) = \kappa_k (I \otimes \sigma_z(k))$ **Error! Filename not specified.**
- **Probabilistic service completion** with rate μ : $A_{\text{dep}} = \mu \sum_n |n-1\rangle\langle n| \otimes \sigma_{\text{dep}}$ **Error! Filename not specified.**

3.4 Critical Enhancements for Quantum Realism

These modeling frameworks can be enhanced to include realistic quantum effects:

- **Fidelity-Weighted Service Rates:** The service rate can be made dependent on a qubit's coherence time, expressed as $\mu \rightarrow \mu \cdot \exp(-t_{\text{wait}}/T_1)$. Here, T_1 is the qubit's coherence time.
- **Entanglement-Aware Metrics:** Success is measured by the probability $P_{\text{succ}} = \text{Tr}[\rho \Pi_{\text{success}}]$ with a projection operator Π_{success} .
- **Resource-Dependent Hamiltonians:** The queue Hamiltonian can be made dependent on available resources, such as $H_{\text{queue}} = \sum_{r \in \text{resources}} g_r(t) H_r$, where $g_r(t)$ allocates qubits and gates.

These frameworks enable the analysis of quantum-specific features like interference in queuing paths, entanglement-induced state correlations, and the fundamental tradeoff between **throughput** $\Lambda = \lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t \langle A_{\text{dep}}^\dagger A_{\text{dep}} \rangle dt$ and **fidelity decay** $F(t) = \exp(-\int_0^t \Gamma(\tau) d\tau)$. For hybrid architectures, these models can be coupled to classical queuing networks via interface operators $B: H_{\text{quant}} \rightarrow X_{\text{class}}$ that measure quantum outputs for classical routing decisions.

4. KEY PERFORMANCE METRICS (QUANTUM-CENTRIC)

Quantum Queueing Theory (QQT) requires a new approach to performance evaluation. Unlike classical systems that rely on simple throughput and latency, QQT uses a dual-dimension framework to quantify both operational efficiency and the integrity of quantum information. These metrics balance temporal efficiency with quantum state preservation.

4.1 Fidelity-Decayed Throughput (Λ_F)

This metric measures the effective task completion rate, weighted by the quality of the output. It accounts for the fact that a completed task is only valuable if its quantum state is preserved.

$$\Lambda_F = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T F_\tau \cdot dN(\tau)$$

Here, $N(\tau)$ is the number of tasks completed by time τ , and $F_\tau = \langle \psi_{\text{ideal}} | \rho_{\text{out}}(\tau) | \psi_{\text{ideal}} \rangle$ is the state fidelity, representing the overlap between the measured output state ρ_{out} and the ideal output state ψ_{ideal} .

4.2 Decoherence-Aware Waiting Time Distribution ($P_{\text{succ}}(t_w)$)

This metric defines the probability of a task succeeding, given its total delay from arrival to service. It directly incorporates the effects of quantum decoherence during the waiting period.

$$P_{succ}(t_w) = P_{circ} \cdot e^{-\Gamma t_w}, \text{ where } \Gamma = k \sum T_1(k)$$

In this equation, t_w is the waiting time, P_{circ} is the inherent success probability of the quantum circuit itself, and $T_1(k)$ are the coherence times of the qubits being used for the task.

4.3 Entanglement Survival Fraction (Ξ)

This metric quantifies how well the quantum correlations (entanglement) between multiple tasks are preserved. It is a critical measure for multi-task applications that require shared entanglement.

$$\Xi = \frac{1}{N} \sum_j \text{Tr}[\rho_j \rho_{j,ent}]$$

Here, ρ_j is the actual output state of task group j , while $\rho_{j,ent}$ is the target entangled state for that group. The trace operation (Tr) measures the overlap between the two states.

4.4 Quantum Resource Utilization (UQ)

Unlike a simple measure of server occupancy, this metric weights server use by the fragility of the quantum resources. It penalizes the underutilization of fragile qubits that have short coherence times.

$$UQ = \int_0^{T_k} \sum_k \langle a_k^\dagger a_k \rangle dt$$

In this formula, a_k and a_k^\dagger are qubit occupancy operators, and $T_1(k)$ is the coherence time for qubit k .

4.5 Stability Condition Augmentation

The classical stability condition is augmented to account for quantum decoherence. This new condition ensures that the system's capacity is sufficient to handle the arrival rate while also maintaining the quality of the quantum states.

$$\lambda_{eff} < \mu \cdot k \cdot \min(e^{-1/(\lambda T_1(k))})$$

Here, λ_{eff} is the effective arrival rate and μ is the service rate. The condition is now limited by the decoherence of the most fragile qubit, with the shortest coherence time $T_1(k)$.

5. HYBRID SYSTEM METRICS

These metrics evaluate the efficiency of systems that combine classical and quantum components, focusing on the interface between the two.

5.1 Classical-Quantum Interface Efficiency (η_{CQ})

This metric measures the quality of quantum state transfer from a quantum server to a classical system relative to the coordination latency.

$$\eta_{CQ} = \text{Avg. classical routing delay} \cdot \text{Tr}[B(\rho_{out})\rho_{class\text{target}}]$$

This formula quantifies the quality of the state transfer and is a crucial measure for a hybrid system's overall performance.

5.2 Quantum Efficiency Frontier

These metrics reveal a fundamental trade-off: maximizing the fidelity-decayed throughput (ΔF) often requires minimizing the waiting time (t_w), which might lead to more errors. This relationship is captured by a new efficiency frontier.

$$\max[\Delta F \cdot P_{\text{succ}}(t_w)] \leq f(\Gamma, \epsilon_g, N_q)$$

The function f represents the upper bound on the product of throughput and success probability, which depends on the decoherence rate (Γ), the gate error rate (ϵ_g), and the number of qubits (N_q). This framework enables the optimization of quantum service systems under the combined constraints of quantum physics and operational demands.

6. CRITICAL RESEARCH DIRECTIONS

6.1. Development of Quantum-Specific Arrival Processes and Service Disciplines

Classical queueing theory relies on well-defined arrival processes (like the Poisson process) and service disciplines (like First-Come, First-Served). QQT needs to formalize quantum analogues that account for superposition and entanglement.

- **Entangled Arrival Streams:** How do we model the arrival of tasks that are entangled with one another? What new queueing behaviors emerge when the arrival of one task is correlated with the state of another? This requires moving beyond classical counting processes to quantum stochastic processes that can handle these non-local correlations.

- **Superpositional Queues and Service:** A task in a superposition of states could require different resources simultaneously or be routed to multiple servers at once. How can a quantum server manage a queue of superpositioned tasks? Research should focus on developing quantum-aware scheduling policies that exploit this parallel nature without causing destructive interference or decoherence. For example, a "Superposition-First, Non-Entangled-Last" (SFNEL) policy could prioritize tasks that are in a superposition state to leverage their inherent parallelism.

6.2. Resource Allocation and Scheduling under Decoherence

Unlike classical servers that are always "on," quantum processors have a finite coherence time. This introduces a new, urgent dimension to resource management.

- **Dynamic, Decoherence-Aware Scheduling:** Standard scheduling algorithms aim to minimize wait times. In QQT, the goal must shift to maximizing the probability of a task's success *before* its qubits decohere. This requires a dynamic scheduling policy that considers not just a task's place in the queue but also its required coherence time (T_1) and gate error rate (ϵ_g). A priority system could be based on a "quantum urgency" score that combines waiting time with the fragility of the qubits required.
- **Optimal Resource Pooling:** Given a limited number of qubits with varying coherence times and fidelities, how should a QPU allocate these resources to incoming tasks? This is a resource optimization problem where the cost is not just a monetary value but the risk of information loss. Research should explore algorithms that dynamically match tasks to the most suitable qubits to maximize overall system fidelity and throughput.

6.3. Stability Analysis of Quantum Queueing Networks

The classical stability condition ($\lambda < \mu$) guarantees that a queue won't grow infinitely large. QQT requires a more nuanced stability analysis that also considers the decay of quantum information.

- **Quantum-Augmented Stability Conditions:** The paper introduces a decoherence-limited stability condition ($\lambda_{\text{eff}} < \mu \cdot \exp(-1/(\lambda T_1))$). This needs to be rigorously derived and expanded for more complex systems, such as multi-server queues and quantum networks. What happens to stability when entanglement and feedback loops are introduced?

- **Stability of Hybrid Quantum-Classical Systems:** In a hybrid architecture, a stable quantum queue doesn't guarantee a stable classical-quantum system. Research is needed to understand how the interface efficiency (η_{CQ}) and classical routing decisions affect the stability of the entire network. This could involve developing a unified theoretical framework that seamlessly merges classical and quantum stability analysis.

6.4. Modeling the Interface Between Quantum and Classical Systems

The most immediate application of QQT will be in hybrid data centers. The transition between quantum and classical processing is a major bottleneck.

- **Quantum-to-Classical State Transfer:** How do we model the process of measuring a quantum state and then using that classical output to make a routing decision in a classical network? This involves creating a formal framework for the **interface operator B** , which is mentioned in the paper but not detailed. This framework should account for measurement overhead and the probabilistic nature of the outcome.
- **Feedback Loops in Hybrid Systems:** Many quantum algorithms (like VQE and QAOA) are iterative and rely on classical feedback. How can QQT models incorporate these continuous quantum-classical feedback loops? This is crucial for accurately predicting the end-to-end performance and latency of such algorithms running on real-world infrastructures.

7. APPLICATIONS & USE CASES

QQT provides a foundational framework for modeling and optimizing the emerging quantum computing ecosystem, with applications ranging from cloud services to specialized quantum networks.

7.1. Quantum Cloud Service Management

QQT is essential for designing and managing **Quantum Computing as a Service (QCaaS)** platforms like those offered by IBM, Amazon, and Microsoft. It helps providers:

- **Predict Performance:** Accurately forecast job completion times, success probabilities, and fidelity for various user workloads, which is crucial for service level agreements (SLAs).
- **Optimize Scheduling:** Implement quantum-aware scheduling algorithms that prioritize tasks based on their qubit requirements, coherence time, and entanglement needs to maximize overall system throughput and job success rates.

- **Pricing Models:** Develop dynamic pricing models that account for the quality and fragility of quantum resources. For instance, a user might pay a premium for a task with a guaranteed high fidelity, which would be prioritized by a QQT-informed scheduler.

7.2. Quantum Data Center and Network Design

As quantum data centers become a reality, QQT will be critical for their physical and operational design.

- **Resource Sizing:** Determine the optimal number of QPUs, cryostats, and support systems required to meet a projected demand, while accounting for the unique constraints of decoherence. This ensures a cost-effective and scalable infrastructure.
- **Network Routing:** Design quantum networks that can route entangled or superpositional tasks between different QPUs. QQT can model the performance of a quantum network where the "hops" are not just a measure of distance but also of fidelity loss due to channel noise.

7.3. Hybrid Quantum-Classical Computing

Many real-world problems require a seamless integration of classical and quantum computation. QQT helps manage the complex interplay between these two environments.

- **Workflow Optimization:** Optimize end-to-end workflows for hybrid algorithms like **Variational Quantum Eigensolver (VQE)** or **Quantum Approximate Optimization Algorithm (QAOA)**. QQT can model the queuing delays and decoherence that occur as data is transferred back and forth between the classical optimizer and the quantum processor.
- **Interface Management:** Design and analyze the performance of the **classical-quantum interface**, which is often a significant bottleneck. QQT can help minimize latency and information loss during the critical process of quantum measurement and classical data transfer.

7.4. Quantum Sensor Networks and Distributed Quantum Systems

QQT can be extended to model distributed systems where quantum information is shared or processed across multiple nodes.

- **Distributed Sensing:** Analyze the performance of a network of quantum sensors where entangled qubits are distributed and must be measured within a specific time

frame to avoid decoherence. QQT metrics like the **Entanglement Survival Fraction (E)** would be key to evaluating system performance.

- **Quantum Internet:** Provide a theoretical basis for the architecture and protocols of a future quantum internet, which will rely on the successful transmission and processing of quantum states. QQT models can help design protocols for managing entangled pairs and minimizing queuing delays in a quantum repeater network.

8. FUTURE DIRECTIONS & OPEN PROBLEMS

8.1. Integrating Error Correction and Fault Tolerance

The current QQT models primarily focus on decoherence as the primary source of error, but they don't fully account for the complexity of quantum error correction (QEC) protocols.

- **Modeling QEC Overhead:** How do you incorporate the significant resource and time overheads of QEC into queueing models? QEC requires a large number of physical qubits to encode a single logical qubit, and these codes have specific latency requirements. Future models must account for this, as it fundamentally changes the service time and resource demands of a task.
- **Impact on Stability and Metrics:** How does the introduction of fault-tolerant quantum computing (FTQC) change the stability conditions and performance metrics? For example, the success probability (P_{succ}) would no longer be a simple exponential decay but would depend on the error-correction code and the rate of logical errors. This would require developing new metrics and stability conditions specifically for FTQC systems.

8.2. Multi-Resource and Networked Queueing Models

The provided models are a good starting point but are limited to a single-server, single-resource perspective. Real-world quantum systems are far more complex.

- **Multi-Qubit and Multi-Server Models:** Develop QQT models that can handle tasks requiring varying numbers of qubits and that can be processed by multiple, interconnected QPUs. This requires formalizing how entanglement can be shared or transferred between different servers, a process known as **entanglement swapping**, and how that affects overall network performance.
- **Quantum Network Routing and Congestion Control:** An open problem is how to route quantum information through a network to minimize decoherence and maximize fidelity. Unlike classical networks, quantum routing is not just about finding the shortest path but the one with the highest fidelity. This calls for the

development of new quantum routing algorithms and congestion control protocols that can handle the unique fragility of quantum states.

8.3. Experimental Validation and Practical Implementation

QQT is currently a theoretical framework. To become a practical tool, its models must be validated against real-world quantum hardware.

- **Simulations and Emulators:** Develop open-source quantum queueing simulators and emulators that can test different scheduling policies and network designs. These tools would need to accurately model real hardware parameters like gate error rates, coherence times, and crosstalk.
- **Collaboration with Quantum Hardware Providers:** A critical step is to partner with companies like IBM, Google, and IonQ to collect real-world queueing data from their quantum computers. This data can be used to validate the assumptions of QQT models and to refine the theoretical framework to better match the realities of current and future quantum hardware. This collaboration is essential for ensuring that QQT remains relevant to the needs of the industry.

8.4. Quantum Machine Learning and Optimization for QQT

Quantum computing itself may provide the tools to solve some of QQT's most complex problems.

- **Quantum-Inspired Optimization:** Explore the use of quantum annealing or hybrid quantum-classical algorithms like QAOA to solve complex resource allocation and scheduling problems that are too difficult for classical computers. For example, a quantum algorithm could find the optimal qubit-to-task assignment to maximize a system's **Quantum Efficiency Frontier** in real time.
- **Learning Quantum Dynamics:** Use machine learning techniques, possibly with quantum data, to learn the complex dynamics of a quantum queueing system. This could help in predicting decoherence rates and other non-Markovian processes that are difficult to model analytically. This approach could lead to more adaptive and robust quantum resource management systems.

9. CONCLUSION

Quantum Queueing Theory (QQT) presents a necessary and timely paradigm shift in service system modeling, addressing the fundamental limitations of classical queueing theory when applied to quantum computing. This paper has formally introduced QQT, outlining a roadmap for future research by defining its core challenges, proposing a suite of

new modeling frameworks, and establishing quantum-centric performance metrics. The unique characteristics of quantum mechanics namely superposition, entanglement, measurement, and decoherence—are not merely minor perturbations but core drivers of a quantum service system's behavior, fundamentally breaking the assumptions of deterministic service times and memoryless processes. By leveraging tools such as quantum stochastic processes, quantum walks, and modified Lindblad master equations, QQT provides a robust theoretical foundation for analyzing and optimizing the performance of a new generation of computing infrastructure. The introduction of metrics like **Fidelity-Decayed Throughput (AF)** and the **Quantum Efficiency Frontier** acknowledges the dual challenge of balancing operational efficiency with the integrity of quantum information. As quantum computing transitions from a scientific curiosity to a practical technology, QQT is poised to become an indispensable tool for engineers, system architects, and researchers. It is the essential framework for designing scalable quantum cloud services, efficient quantum data centers, and the integrated hybrid quantum-classical networks that will power the next era of computation. The path forward involves tackling open problems in quantum error correction, developing multi-resource network models, and validating theoretical predictions against real-world hardware, thereby solidifying QQT as a critical discipline for the future of information technology.

REFERENCES:

1. Nielsen, M. A., & Chuang, I. L. (2010). *Quantum computation and quantum information* (10th ed.). Cambridge University Press.
2. Preskill, J. (2018). Quantum computing in the NISQ era and beyond. *Quantum*, *2*, 79.
3. Shor, P. W. (1997). Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. *SIAM Journal on Computing*, *26*(5), 1484–1509.
4. Harrow, A. W., Hassidim, A., & Lloyd, S. (2009). Quantum algorithm for linear systems of equations. *Physical Review Letters*, *103*(15), 150502.
5. Peruzzo, A., et al. (2014). A variational eigenvalue solver on a photonic quantum processor. *Nature Communications*, *5*, 4213.
6. Farhi, E., Goldstone, J., & Gutmann, S. (2014). A quantum approximate optimization algorithm. *arXiv:1411.4028*.
7. Terhal, B. M. (2015). Quantum error correction for quantum memories. *Reviews of Modern Physics*, *87*(2), 307–346.
8. Campbell, E. T., Terhal, B. M., & Vuillot, C. (2017). Roads towards fault-tolerant universal quantum computation. *Nature*, *549*(7671), 172–179.
9. Gottesman, D. (2010). An introduction to quantum error correction and fault-tolerant quantum computation. *Proceedings of Symposia in Applied Mathematics*, *68*, 13–58.
10. de Leon, N. P., et al. (2021). Materials challenges for quantum technologies based on solid-state defects. *Nature Reviews Materials*, *6*(10), 906–925.
11. Wehner, S., Elkouss, D., & Hanson, R. (2018). Quantum internet: A vision for the road ahead. *Science*, *362*(6412), eaam9288.

12. Kimble, H. J. (2008). The quantum internet. *Nature*, *453*(7198), 1023–1030.
13. Chakraborty, K., Rozpedek, F., Dahlberg, A., & Wehner, S. (2019). Distributed routing in a quantum internet. *arXiv:1907.11630*.
14. Gyongyosi, L., & Imre, S. (2019). A survey on quantum computing technology. *Computer Science Review*, *31*, 51–71.
15. Kleinrock, L. (1975). *Queueing systems, volume 1: Theory*. Wiley.
16. Harchol-Balter, M. (2013). *Performance modeling and design of computer systems: Queueing theory in action*. Cambridge University Press.
17. Botea, A., Bradu, A., & Saffidine, A. (2021). Scheduling for quantum computers. *Proceedings of the International Conference on Automated Planning and Scheduling*, *31*, 55–63.
18. Li, Z., et al. (2020). Quantum scheduling for critical tasks in distributed quantum computing. *IEEE Transactions on Quantum Engineering*, *1*, 1–12.
19. Breuer, H. P., & Petruccione, F. (2007). *The theory of open quantum systems*. Oxford University Press.
20. Gardiner, C., & Zoller, P. (2004). *Quantum noise* (3rd ed.). Springer.
21. Venegas-Andraca, S. E. (2012). Quantum walks: A comprehensive review. *Quantum Information Processing*, *11*(5), 1015–1106.
22. Portugal, R. (2018). *Quantum walks and search algorithms*. Springer.
23. Kretschmer, W., & Luongo, G. (2021). Quantum queueing. *arXiv:2103.12076*.
24. Arute, F., et al. (2019). Quantum supremacy using a programmable superconducting processor. *Nature*, *574*(7779), 505–510.
25. Gambetta, J. M., Chow, J. M., & Steffen, M. (2017). Building logical qubits in a superconducting quantum computing system. *npj Quantum Information*, *3*(1), 2.
26. Cross, A. W., Bishop, L. S., Smolin, J. A., & Gambetta, J. M. (2017). Open quantum assembly language. *arXiv:1707.03429*.
27. Smith, R. S., Curtis, M. J., & Zeng, W. J. (2016). A practical quantum instruction set architecture. *arXiv:1608.03355* (Rigetti).
28. Proctor, T., et al. (2021). Scalable randomized benchmarking of quantum computers using mirror circuits. *arXiv:2112.09853*.
29. Cerezo, M., et al. (2021). Variational quantum algorithms. *Nature Reviews Physics*, *3*(9), 625–644.
30. Endo, S., Cai, Z., Benjamin, S. C., & Yuan, X. (2021). Hybrid quantum-classical algorithms and quantum error mitigation. *Journal of the Physical Society of Japan*, *90*(3), 032001.
31. Kandala, A., et al. (2017). Hardware-efficient variational quantum eigensolver for small molecules and quantum magnets. *Nature*, *549*(7671), 242–246.
32. Bharti, K., et al. (2022). Noisy intermediate-scale quantum algorithms. *Reviews of Modern Physics*, *94*(1), 015004
33. Biamonte, J., et al. (2017). Quantum machine learning. *Nature*, *549*(7671), 195–202.
34. Havlíček, V., et al. (2019). Supervised learning with quantum-enhanced feature spaces. *Nature*, *567*(7747), 209–212.
35. Rebentrost, P., Mohseni, M., & Lloyd, S. (2014). Quantum support vector machine for big data classification. *Physical Review Letters*, *113*(13), 130503.

36. Montanaro, A. (2016). Quantum algorithms: An overview. *npj Quantum Information*, *2*(1), 1–8.
37. Yalin, G., & Kais, S. (2021). Resource management for quantum computing. *arXiv:2105.08202*.
38. Das, S., & Chakrabarti, B. K. (2008). Quantum annealing and related optimization methods. *Lecture Notes in Physics*, *679*, 3–38.
39. Van Meter, R., & Horsman, C. (2013). A blueprint for building a quantum computer. *Communications of the ACM*, *56*(10), 84–93.
40. Córcoles, A. D., et al. (2021). Challenges and opportunities of near-term quantum computing systems. *Proceedings of the IEEE*, *108*(8), 1338–1352.
41. Aspuru-Guzik, A., Dutoi, A. D., Love, P. J., & Head-Gordon, M. (2005). Simulated quantum computation of molecular energies. *Science*, *309*(5741), 1704–1707.
42. Cao, Y., et al. (2019). Quantum chemistry in the age of quantum computing. *Chemical Reviews*, *119*(19), 10856–10915.
43. Wilde, M. M. (2017). *Quantum information theory* (2nd ed.). Cambridge University Press.
44. Acín, A., et al. (2018). The quantum technologies roadmap: A European community view. *New Journal of Physics*, *20*(8), 080201.
45. Devitt, S. J., Munro, W. J., & Nemoto, K. (2013). Quantum error correction for beginners. *Reports on Progress in Physics*, *76*(7), 076001.
46. Zhang, J., & Zhang, K. (2019). Queueing analysis for quantum communications. *IEEE Access*, *7*, 152833–152846.
47. IBM Quantum Experience. (2021). *IBM Quantum computing platform*. <https://quantum-computing.ibm.com/>
48. Amazon Braket. (2020). *Amazon Web Services*. <https://aws.amazon.com/braket/>
49. Abhari, A. J., et al. (2012). Scaffold: Quantum programming language. *Technical Report*, Princeton University.