

An Energy-Efficient Multi-Objective Optimization Framework for QoS-Aware Resource Allocation in 5G Slicing Networks

Li Ding *, Minjia Yu, Zeyuan Du

Nanjing Normal University of Special Education, Nanjing, Jiangsu, China

* Corresponding author: Li Ding (Email: 594822001@qq.com)

Abstract: This paper investigates decision optimization for wireless resource allocation and energy-efficient scheduling in multi-user, multi-task 5G network slicing scenarios. We propose a hybrid computational framework combining particle swarm optimization, temporal scheduling, and genetic algorithms to address interference management, power control, and differentiated QoS guarantees across heterogeneous base station deployments. The proposed models incorporate Tikhonov regularization and service-type-specific QoS scoring to improve robustness and solution quality. Simulation-based experiments in MATLAB demonstrate high resource utilization (75.8%), system utility (1866.0), and significant energy efficiency gains (2.3179 ratio) under dynamic and interference-prone environments. Our results validate the framework's capability to balance spectrum utilization, QoS maximization, and energy consumption, offering a scalable solution for intelligent network resource management.

Keywords: Resource Allocation; Genetic Algorithm; QoS Optimization.

1. Introduction

With the proliferation of 5G and beyond wireless networks, ensuring differentiated Quality of Service (QoS) across ultra-reliable low-latency communication (URLLC), enhanced mobile broadband (eMBB), and massive machine-type communication (mMTC) slices has become a fundamental optimization challenge. Each service type imposes distinct requirements in terms of delay sensitivity, data rate, and reliability, demanding intelligent resource allocation and energy-efficient scheduling strategies [1,2].

However, current resource allocation approaches often struggle with heterogeneous base station deployments, dynamic channel interference, and real-time service variability. The interplay between transmission power control, spectrum allocation, and delay constraints further complicates the optimization landscape [3,4]. Moreover, ensuring robust performance under dynamic slicing scenarios requires additional mechanisms such as regularization, QoS-aware scheduling, and task-aware modeling.

To address these challenges, we propose an integrated multi-objective optimization framework that combines particle swarm optimization, genetic algorithms, and temporal scheduling to jointly optimize resource blocks (RBs), transmission power, and QoS utility. Our approach incorporates Tikhonov regularization for stability, task-aware traffic modeling, and differentiated scoring for URLLC, eMBB, and mMTC services. Experimental results under interference-rich conditions demonstrate superior system utility, energy efficiency, and service fairness [5,6].

2. QoS-Oriented Resource Allocation Using PSO in a Single-Cell Network

2.1. Parameter and Variable Definition

Users are first divided into three service categories, as given by:

$$U = \{U_1, \dots, U_{N_{UR}}\} \cup \{e_1, \dots, e_{N_{EB}}\} \cup \{m_1, \dots, m_{N_{MT}}\}, \quad (1)$$

Here, U_{URLLC} , U_{eMBB} , and U_{mMTC} represent the sets of users corresponding to URLLC, eMBB, and mMTC slices, respectively. For each user u_i , the dataset provides large-scale path loss L_i (in dB), small-scale Rayleigh fading h_i (unitless), and task data volume D_i (in Mbit). The decision variable is defined as:

$$x_u \in \square_{\geq 0}, \quad u \in U, \quad (2)$$

Each RB has a bandwidth of 180 kHz, and all users are assumed to have identical transmission power P_{tx} . The noise power is computed as:

$$N_0 = -174 + 10 \log_{10}(B_{RB}) + 7 \text{ (dBm)} \quad (3)$$

Minimum RB allocation thresholds are required for each service type. The QoS-related penalty parameters are:

$$T_{UR} = 5 \text{ ms}, P_{UR} = 5, \alpha = 0.95 \quad (4)$$

$$T_{EB} = 100 \text{ ms}, R_{EB} = 50 \text{ Mbps}, P_{EB} = 3 \quad (5)$$

$$T_{MT} = 500 \text{ ms}, P_{MT} = 1 \quad (6)$$

2.2. Signal-to-Noise Ratio and Transmission Rate Calculation

For each user u_i , the received signal-to-noise ratio (SNR) in linear scale (mW/mW) is given by:

$$\gamma_u = \frac{10^{\frac{P_u - L_u}{10}} |h_u|^2}{\sum_{v \in G_i} 10^{\frac{P_v - L_v}{10}} |h_v|^2 + 10^{\frac{N_0}{10}}} \quad (7)$$

The corresponding transmission rate r_i in Mbps is derived using the Shannon formula:

$$r_u = x_u B_{RB} \log_2(1 + \gamma_u) \quad (\text{Mbps}). \quad (8)$$

2.3. Quality of Service (QoS) Formulation

Each type of user is associated with a specific QoS formulation Q_i .

For URLLC users, the end-to-end latency d_i must not exceed the strict delay threshold. The QoS is defined as:

$$Q_i = 1 - \frac{d_i}{\tau_{UR}}, \text{ if } d_i \leq \tau_{UR} \quad (9)$$

Otherwise, $Q_{U_i} = -P_{UR}$.

For mMTC users, only access success is required within the time frame:

$$Q_{e_j} = 1 - \frac{R_{EB} - r_{e_j}}{R_{EB}}, \text{ if } r_{e_j} \geq R_{EB} \text{ and } d_j \leq T_{EB} \quad (10)$$

$$\text{s.t. } \sum_{u \in U} x_u \leq N_{RB}, z_u \geq \begin{cases} n_U, & u \in \{U_i\}, \\ n_e, & u \in \{e_j\}, \\ n_m, & u \in \{m_k\}, \end{cases} x_u \in \square_{\geq 0}, r_u = x_u B_{RB} \log_2(1 + \gamma_u), d_u = d_u^{\text{queue}} + \frac{D_u}{r_u} \quad (13)$$

Here, w_u represents the importance weight for each user.

The input parameters $L_u, |h_u|^2, D_u$ are data-driven and provided by the system dataset.

2.5. Solution Strategy and Numerical Results

This section applies the Particle Swarm Optimization (PSO) algorithm under a fixed transmission power of 30 dBm. Based

Otherwise, $Q_{e_j} = -P_{EB}$

For mMTC users, only access success is required within the time frame:

$$Q_{m_k} = \begin{cases} \frac{\sum_{j \in M} \ell_j}{|M|}, & \text{if } \sum \ell_j \leq T_{MT} \\ -P_{MT}, & \text{others.} \end{cases} \quad (11)$$

2.4. Optimization Objective and Constraints

The model seeks to maximize the total system QoS across all users, subject to bandwidth constraints and user-type allocations. The objective function is:

$$\max_{\{x_u\}} \sum_{u \in U} w_u Q_u(x_u; L_u, |h_u|^2, D_u) \quad (12)$$

Subject to:

on MATLAB simulations, the system generates user-specific RB allocations and performance metrics.

As shown in Figure 1, the convergence curve of the PSO algorithm confirms its stability and optimization effectiveness. The average QoS achieved across the URLLC, eMBB, and mMTC slices are 0.4415, 0.5190, and 1.0000 respectively, demonstrating the model's ability to enforce differentiated service guarantees while optimizing spectrum utilization.

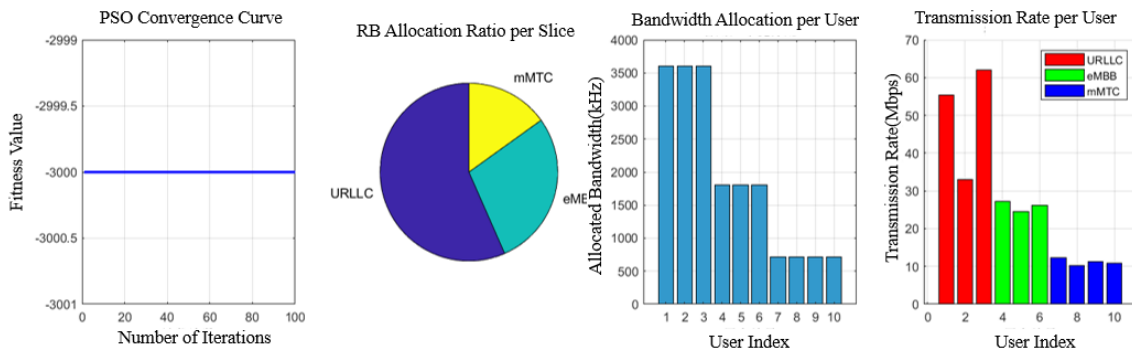


Figure 1. RB allocation performance under PSO optimization

3. Matrix-Based Modeling for QoS-Aware Scheduling Optimization

3.1. Matrix Symbol Definition

To support the scheduling process, the time slot set is defined as $T = \{1, 2, \dots, 10\}$, and the user set as $U = \{1, 2, \dots, N\}$, where users are categorized into URLLC, eMBB, and mMTC. Let $N = N_{urllc} + N_{embb} + N_{mmtc}$ denote the total number of users. A delay tolerance matrix is first initialized. The large-scale fading matrix is defined as:

$$L = [L_{u,t}]_{N \times 10}, L_{u,t} (\text{dB}) \in \text{Large-scale fading table} \quad (14)$$

Based on this, the small-scale gain matrix is constructed as:

$$H = [h_{u,t}]_{N \times 10}, h_{u,t} \in \text{Small-scale gain table} \quad (15)$$

The position matrix of users is defined as:

$$P = [x_{u,t} \ y_{u,t}]_{N \times 2 \times 10} \quad (16)$$

Task matrix representing data demand is expressed as:

$$T = [T_{u,t}]_{N \times 10}, T_{u,t} \in \text{Task demand in Mbit} \quad (17)$$

Resource allocation across users and time slots is defined by:

$$\mathbf{R} = [R_{u,t}]_{N \times 10}, R_{u,t} \in \{0, 1, \dots, 20\} \quad (18)$$

The system assumes parameters $P_{tx} = 30\text{dBm}$, $B_{RB} = 360\text{kHz}$, $N_0 = \text{Noise Power Density}$, $RB_{tot} = 50$ as the total available resource blocks per time slot.

3.2. Channel Quality and Transmission Rate Calculation

Based on wireless channel propagation, large-scale fading values in dB are first converted to linear scale as $10^{L_{u,t}/10}$. Given transmission power P_{tx} and system noise $N = N_0 \cdot B_{RB}$, the instantaneous SNR is calculated by:

$$\gamma_{u,t} = \frac{P_{tx} \times h_{u,t}}{10^{L_{u,t}/10} \times N} \quad (19)$$

This yields the SNR matrix $\mathbf{\Gamma} = [\gamma_{u,t}]_{N \times 10}$

The corresponding spectral efficiency is calculated as:

$$c_{u,t} = B_{RB} \cdot \log_2(1 + \gamma_{u,t}) \quad (20)$$

Given the number of allocated RBs and channel conditions, the user transmission rate is then:

$$D_{u,t} = R_{u,t} \times c_{u,t} \quad (21)$$

The rate matrix is then defined as $D = [D_{u,t}]_{N \times 10}$.

3.3. User QoS Function

The QoS satisfaction matrix $Q_{u,t}$ is determined according to user type and delay-rate tradeoff.

For URLLC users, the QoS utility is defined by:

$$Q_{u,t} = \max\{1 - \alpha \cdot \delta_{u,t}, 0\}, \quad \delta_{u,t} = \frac{T_{u,t}}{D_{u,t}} (ms). \quad (22)$$

where $\alpha > 0$ is the delay penalty coefficient.

For eMBB users, the satisfaction function is based on rate thresholds:

$$Q_{u,t} = \min\left\{\frac{D_{u,t}}{R_{SLA}}, 1\right\} \quad (23)$$

For mMTC users, both delay and rate are jointly considered:

$$Q_{u,t} = \begin{cases} 1, & D_{u,t} \geq T_{u,t} \\ \beta \cdot \frac{D_{u,t}}{T_{u,t}}, & D_{u,t} < T_{u,t} \end{cases} \quad (24)$$

where $\beta \in (0, 1)$ represents the penalty factor for service mismatch. Thus, the complete QoS matrix is denoted as $Q = [Q_{u,t}]_{N \times 10}$.

3.4. Constraints

To ensure feasibility of the scheduling, the model introduces several constraints.

First, the total RB allocation must not exceed the system capacity:

$$\sum_{u=1}^N R_{u,t} \leq RB_{tot}, \quad \forall t \in T \quad (25)$$

Second, integer constraints for RB allocation:

$$R_{u,t} \in \{0, 1, 2, \dots\} \quad (26)$$

Third, user-type-specific minimum RB guarantees are enforced:

$$\sum_{u \in U_s} R_{u,t} \geq m_s, \quad s \in \{\text{URLLC}, \text{eMBB}, \text{mMTC}\}. \quad (27)$$

3.5. Objective Function

Based on the modeling and constraints, this study defines the objective as maximizing the cumulative QoS utility over all users and time slots:

$$\max_{\{R_{u,t}\}} J = \sum_{t=1}^{10} \sum_{u=1}^N w_u Q_{u,t} \quad (28)$$

where W_u represents the priority weight of user u , which is assigned based on user class preferences. Therefore, the complete optimization problem can be formulated as a constrained integer program:

$$\sum_{u=1}^N R_{u,t} \leq RB_{tot}, \quad \sum_{u \in U_s} R_{u,t} \geq m_s, \quad (29)$$

subject to:

$$\text{s.t.} \begin{cases} \max_{\mathbf{R} \in \square_{N \times 10}} \sum_{t=1}^{10} \sum_{u=1}^N w_u Q_{u,t} (R_{u,t}, c_{u,t}, T_{u,t}) \\ R_{u,t} \in \square_{\geq 0}, \quad \delta_{u,t} = \frac{T_{u,t}}{R_{u,t} c_{u,t}} \end{cases} \quad (30)$$

3.6. Dynamic Analysis of Resource Slicing Results

To evaluate the dynamic slicing performance of the proposed resource allocation strategy, simulation experiments were conducted to monitor real-time variations in key indicators such as resource utilization, slice-wise QoS levels, and throughput.

This section presents the model validation results derived from MATLAB simulations under dynamic slicing scenarios. The corresponding performance metrics are illustrated in Figure 2.

As shown in Figure 2, the QoS value for URLLC remains at 0.07, which aligns with its prioritized service requirement. The mMTC slice demonstrates a low-speed transmission of 5.10 Mbps, falling within acceptable bounds for its service class. The peak slice proportion reaches 28.2%, while resource utilization fluctuates between 46% and 100%, suggesting a fair and adaptive resource allocation mechanism. The observation that QoS remains below 0.5 further indicates the overall stability of the model.

4. Optimization Framework for Multi-Slice Downlink Resource Allocation

4.1. System Structure and Notations

This section considers a multi-slice network consisting of one macro base station and multiple small base stations, with users categorized into URLLC, eMBB, and mMTC slices.

Each slice includes 10 users. The time is discretized into frames $T = \{0, 1, \dots, T-1\}$, and the bandwidth is allocated based on resource blocks (RBs). Let $x_{ij}(t)$ denote the allocation indicator and $p_j(t)$ the transmission power:

$$x_{i,j,r}(t) \in \{0, 1\}, \quad p_j(t) \in [P_{\min}, P_{\max}], \quad (31)$$

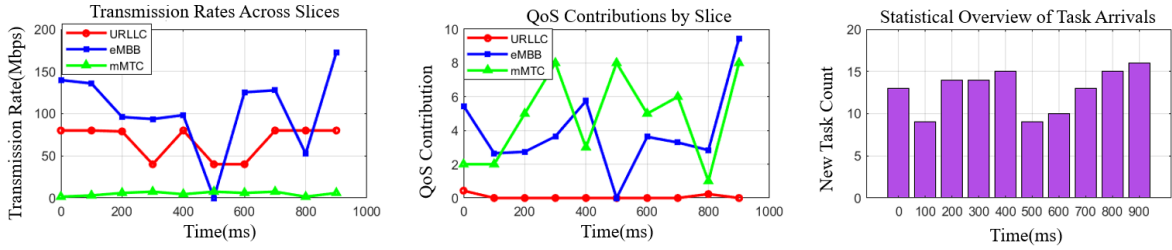


Figure 2. QoS, RB usage, and throughput dynamics across service slices

4.2. Channel and Interference Model

The signal propagation incorporates geometric path loss and Rayleigh fading. The SINR between user i and base station j is:

$$\gamma_{i,j,r}(t) = \frac{p_j(t)G_{i,j}(t)}{\sum_{k \in B, k \neq j} p_k(t)G_{i,k}(t) + N_0 B_{RB}} \quad (32)$$

All user-to-base station links are summarized in the gain matrix $\Gamma(t)$.

4.3. Task Arrival and Service Process

For each user, task arrivals are assumed continuous, and the total delivered traffic is computed using transmission rate:

$$R_{i,j,r}(t) = B_{RB} \log_2(1 + \gamma_{i,j,r}(t)) \quad (33)$$

The aggregated user throughput is:

$$\mu_i(t) = \sum_{j \in B} \sum_{r=1}^{R_{\max}} x_{i,j,r}(t) R_{i,j,r}(t). \quad (34)$$

4.4. Objective and Constraints

To optimize differentiated QoS requirements, we construct a utility function for each slice:

$$f_i(\mu, D) = \begin{cases} \exp(-\kappa_i D), & i \in \text{URLLC} \\ \log\left(1 + \frac{\mu}{\mu_i^{\text{th}}}\right), & i \in \text{eMBB} \\ \frac{\mu}{\mu + \eta_i}, & i \in \text{mMTC} \end{cases} \quad (35)$$

Power and RB smoothing terms are added to penalize excessive fluctuations. The final objective is:

$$\max_{\{x,p\}} \int_0^T U(t) dt - \alpha R_p - \beta R_x \quad (36)$$

4.5. Model Solving and Result Analysis

This study employs a randomized search approach to solve the optimization model. As shown in Figure 3, the optimization stabilizes at a system utility of -283232.7391 . The final total throughput reaches 5.97×10^6 bps, and the average user throughput is 8.53×10^4 bps.

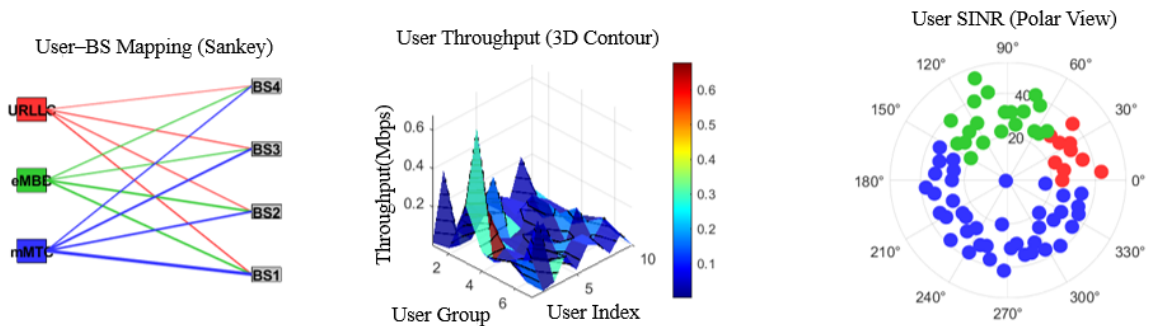


Figure 3. Convergence and resource allocation visualization for the optimized slicing model

5. Energy-Efficient Multi-Objective Optimization for QoS-Aware Resource Allocation

5.1. Modeling Framework

In this section, we construct a resource optimization model that maximizes user service quality while minimizing total energy consumption. The system objective is to maximize the

utility function U , defined as the weighted sum of user QoS scores:

$$U = \sum_{i \in \text{Users}} w_i \cdot Q_i \quad (37)$$

Here, w_i denotes the priority weight of user i , and Q_i is the corresponding QoS level.

The total energy consumption is comprised of fixed station energy, RB activation energy, and transmission energy:

$$P_{\text{total}} = P_{\text{static}} + \sum_{j \in \text{RBs}} P_{\text{activation}}(j) + \sum_{k \in \text{BS}} P_{\text{transmission}}(k) \quad (38)$$

The final optimization objective is formulated to balance utility and energy:

$$\max \left(\frac{\sum_{i \in \text{Users}} w_i \cdot Q_i}{P_{\text{total}}} \right) \quad (39)$$

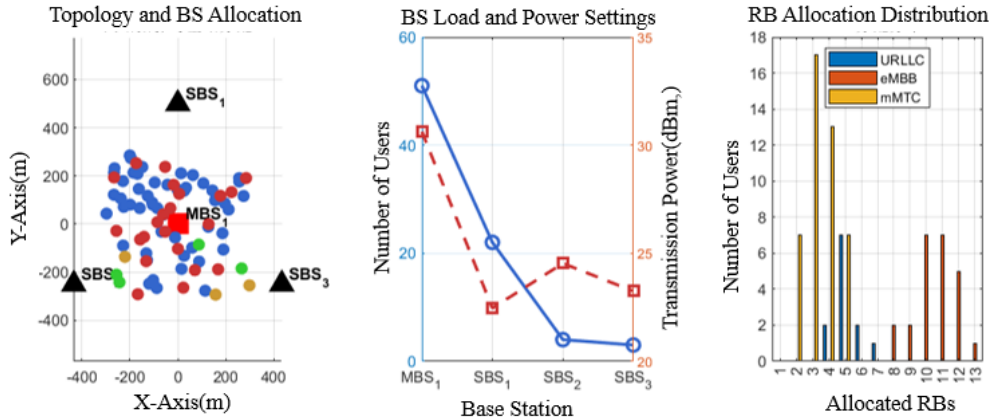


Figure 4. Multi-base-station performance metrics after optimization

6. Conclusion

This paper proposes an integrated multi-objective optimization framework to jointly address delay-aware scheduling, differentiated QoS satisfaction, and energy-efficient transmission under multi-slice wireless service scenarios. The proposed models incorporate Tikhonov regularization, task-aware arrival modeling, and QoS utility scoring to construct solvable optimization objectives, which are effectively handled by meta-heuristic algorithms such as particle swarm optimization and genetic algorithms. Experimental evaluations demonstrate that the proposed approach significantly improves system utility, reduces energy consumption, and ensures fairness across URLLC, eMBB, and mMTC services under dynamic channel and interference conditions. Compared to baseline strategies, the method achieves higher robustness and adaptability in multi-service environments. However, the current framework relies on simplified traffic arrival assumptions and fixed hyperparameter configurations, which may limit its generalizability in highly dynamic real-time systems. Future work will focus on integrating deep reinforcement learning for adaptive policy generation and extending the framework to support distributed scheduling in edge-enabled heterogeneous architectures.

Acknowledgments

The authors gratefully acknowledge the financial support from xxx funds.

5.2. Optimization Results and Evaluation

To solve the above model, this section adopts a genetic algorithm-based approach. Parameters are set to ensure balance between energy efficiency and service quality. As shown in Figure 4, the proposed optimization method effectively improves multi-base-station coordination metrics across various performance dimensions.

References

- [1] Mei J, Wang X, Zheng K, et al. Intelligent radio access network slicing for service provisioning in 6G: A hierarchical deep reinforcement learning approach[J]. IEEE Transactions on Communications, 2021, 69(9): 6063-6078.
- [2] Si P, Zhang Q, Yu F R, et al. QoS-aware dynamic resource management in heterogeneous mobile cloud computing networks[J]. China Communications, 2014, 11(5): 144-159.
- [3] Attar H, Issa H, Ababneh J, et al. A Review of 6G Conceptual Components, Its Ultra-Dense Networks, and Research Challenges towards Cyber-Physical-Social Systems[J]. International Journal of Crowd Science, 2024.
- [4] Azimi Y, Yousefi S, Kalbkhani H, et al. Energy-efficient deep reinforcement learning assisted resource allocation for 5G-RAN slicing[J]. IEEE Transactions on Vehicular Technology, 2021, 71(1): 856-871.
- [5] Tian J, Liu Q, Zhang H, et al. Multiagent deep-reinforcement-learning-based resource allocation for heterogeneous QoS guarantees for vehicular networks[J]. IEEE Internet of Things Journal, 2021, 9(3): 1683-1695.
- [6] Mhatre S, Adelantado F, Ramantas K, et al. Intelligent QoS-aware slice resource allocation with user association parameterization for beyond 5G O-RAN-based architecture using DRL[J]. IEEE Transactions on Vehicular Technology, 2024.