Enhancing 4G/5G Networks: Instantaneous Uplink Estimation and Scalable RAN Slicing Strategies

[1] Ali Zahir Yassin Arabi [2] Jitendra Vaswani

[1] Post Graduate Scholar [2] Assistant Professor

Department of ECE (Electronics and Communication Engineering), Mewar University, Chittorgarh, Rajasthan, India

Abstract

The rapid evolution of cellular networks, particularly 4G and 5G, has been driven by the demands of diverse sectors requiring support for sophisticated services and high-volume data traffic. However, these networks face significant challenges in meeting stringent performance expectations while handling massive uplink and downlink loads. This thesis addresses these challenges through two key contributions aimed at enhancing the intelligence and flexibility of cellular networks. First, it introduces an intelligent framework for estimating users' instantaneous uplink throughput at fine-grained time intervals. A scalable estimation model leveraging machine learning techniques-including Linear Regression, Random Forest, and Support Vector Regression—is developed and validated on data gathered from a real-time 4G testbed simulating diverse radio conditions. Results indicate high estimation accuracy, with errors under 15%, particularly for forecast windows exceeding 700 ms, while highlighting the insufficiency of radio measurements alone for precise predictions at smaller timescales. The second contribution focuses on enforcing 5G Radio Access Network (RAN) slicing at the resource level from a multi-cell perspective. While core network slicing benefits from cloud-based solutions, RAN slicing faces complex challenges related to slice orthogonality, satisfaction, scalability, and cooperation. An exact optimization model based on constraint programming is proposed, alongside a 2D bin packing heuristic, to balance these competing requirements. Additionally, three heuristics are introduced to prioritize scalability without sacrificing key performance metrics. Experimental results demonstrate strong performance, particularly with two heuristics, underscoring their effectiveness in enabling real-time RAN slicing. Collectively, these contributions advance the capabilities of 4G/5G networks to meet modern service demands.

Keywords: 4G/5G Networks, Machine Learning, Uplink Throughput Estimation, RAN Slicing, Resource Optimization

1. Introduction

The exponential growth in data-centric applications and the emergence of latency-sensitive, highthroughput services have placed unprecedented demands on the performance, reliability, and adaptability of 4G and 5G cellular networks. The uplink path, in particular, is becoming increasingly critical with the proliferation of machine-type communication (MTC), industrial IoT, real-time video analytics, and mission-critical services. These paradigms necessitate fine-grained, accurate estimation of uplink throughput to enable proactive radio resource management, meet Quality of Service (QoS) guarantees, and support advanced scheduling algorithms.

This research investigates two fundamental challenges in modern cellular network optimisation: (i) real-time estimation of users' instantaneous uplink throughput using machine learning (ML) techniques, and (ii) scalable and constraint-aware radio access network (RAN) slicing. To address the first, we design a predictive framework that leverages a live 4G testbed to collect diverse lower-layer (PHY/MAC) eNB metrics under realistic channel and traffic conditions. The dataset is used to train and evaluate several supervised ML models—namely Linear Regression (LR), Random Forest (RF), and Support Vector Regression (SVR)—to forecast uplink throughput at millisecond-level resolution. Empirical results show that ML models achieve estimation errors below 15% for time windows above 700 ms, though estimation accuracy degrades significantly at shorter granularities due to intrinsic radio variability and limitations of the available signal-level features.

The second contribution addresses the increasingly critical need for RAN slicing in 5G networks, where multiple virtualised services with heterogeneous requirements must coexist over a shared physical infrastructure. Unlike core network slicing, RAN slicing enforcement is constrained by inter-slice orthogonality, SLA compliance, scalability, and inter-cell cooperation. We formulate a constraint programming (CP)-based exact model that guarantees optimal resource allocations under all four constraints. Additionally, a scalable 2D bin-packing heuristic is introduced, optimising for slice satisfaction, orthogonality, and scalability, albeit without explicit cooperation mechanisms. To bridge this gap, three novel heuristics are proposed, enabling simultaneous enforcement of all slicing constraints. Experimental validation confirms that two of the proposed heuristics offer effective trade-offs between computational efficiency and solution quality, making them viable for near-real-time RAN slicing scenarios.

Collectively, this work contributes an integrated approach to enhancing the intelligence and flexibility of cellular networks. By combining machine learning-based throughput forecasting with constraint-aware RAN slicing mechanisms, the study provides key enablers for the next generation of network-aware services and adaptive resource orchestration in 4G and 5G environments.

The evolution of mobile telecommunications has been marked by significant generational shifts, transforming how societies communicate, exchange data, and enable new digital economies. Fourth Generation (4G) networks, widely deployed globally, have provided substantial

enhancements in mobile broadband performance, supporting high data rates, low latency, and a range of services that fuelled the proliferation of smartphones and rich multimedia applications [1]–[4]. However, with the relentless growth in user demand, emerging applications such as the Internet of Things (IoT), vehicular networks, and Industry 4.0, and the ambition for ubiquitous connectivity, the limitations of 4G systems have become increasingly evident [5]–[8].

Fifth Generation (5G) networks have emerged to address these challenges, promising not merely incremental improvements but transformative capabilities. While 4G largely focused on enhancing human-centric communication services, 5G expands the horizon towards machine-type communications, ultra-reliable low-latency communication (URLLC), and massive connectivity [9]–[12]. This paradigm shift underpins diverse applications ranging from autonomous vehicles [27], smart manufacturing [24], to augmented and virtual reality [4], demanding stringent performance parameters in throughput, latency, and reliability.

Technologically, 5G introduces significant innovations over its predecessor, including utilisation of millimetre-wave (mmWave) frequencies [4], advanced multiple-input multiple-output (MIMO) schemes [17], and novel access techniques such as non-orthogonal multiple access (NOMA) [9], [26]. These enable higher spectral efficiency and the capacity to support vastly greater numbers of devices concurrently [15], [19]. However, these advancements come with considerable challenges in terms of network design, energy consumption [23], security [13], and efficient spectrum management [6].

One of the key architectural developments in 5G is the concept of network slicing, allowing operators to create virtualised, end-to-end networks tailored for specific use cases, thereby improving resource utilisation and service flexibility [12], [16]. Moreover, edge computing has become integral to 5G, pushing computation closer to end-users to reduce latency and enhance quality of experience for applications like real-time video analytics and autonomous systems [8], [16].

Researchers have extensively explored the trade-offs between energy efficiency and performance in 5G deployments, recognising sustainability as a critical design objective [17], [23]. Similarly, significant attention has been devoted to developing robust security frameworks, given the expanded attack surface introduced by massive device connectivity and virtualised architectures [13], [27].

Despite the significant strides achieved with 5G, the research community continues to identify areas necessitating further innovation, including seamless integration with existing 4G infrastructures, optimisation of machine learning techniques for network management [14], [22], and preparation for the eventual evolution towards Sixth Generation (6G) networks [3], [21]. Performance comparisons have illustrated the considerable gains 5G offers over 4G, yet practical deployment scenarios continue to reveal gaps between theoretical promises and real-world performance, particularly in coverage at higher frequencies and cost-effectiveness [4], [29].

In summary, while 4G networks have laid the essential groundwork for mobile broadband and digital transformation, 5G represents a crucial leap forward, both technologically and architecturally, to fulfil the ambitious requirements of modern and future applications [1]–[30]. It is within this evolving landscape that research into the comparative performance, capabilities, and challenges of 4G and 5G networks remains highly relevant and critical for guiding ongoing development and deployment strategies.

2. Methodology

In this work, we propose a methodology for estimating uplink throughput in 4G networks based on eNB lower-layer metrics using supervised machine learning techniques (MLTs), namely Linear Regression (LR), Support Vector Regression (SVR), and Random Forest (RF), selected to capture both linear and non-linear relationships between network metrics and throughput. LR models linear dependencies, SVR leverages kernel functions such as the Radial Basis Function (RBF) to handle non-linearities while optimising parameters like ε , γ , and C for balanced bias-variance trade-offs, and RF aggregates predictions from multiple decision trees to improve robustness against over-fitting, with key hyper-parameters including the number of estimators and tree depth. To identify optimal hyper-parameters efficiently, we employ random search combined with Kfold cross-validation (CV), where K=10 is adopted for reliable error estimation, and further mitigate estimation bias through nested CV, which uses an inner loop for hyper-parameter tuning and an outer loop for generalisation error assessment. Additionally, we introduce two temporal modelling parameters: a forecast window (β), enabling throughput predictions over variable future intervals, and a lag window (l), capturing historical metric patterns to enhance prediction accuracy. For each MLT and combination of β and l, the methodology applies nested CV to determine the model with the lowest root mean square error (RMSE), thereby ensuring accurate, generalisable throughput estimation across diverse network conditions.

3. Instantaneous Uplink Throughput Estimation (Technical Single Paragraph)

Instantaneous uplink throughput estimation seeks to predict the achievable data rate of a user equipment (UE) within a fine-grained temporal interval, leveraging evolved Node B (eNB) lower-layer metrics that capture the physical and MAC layer conditions of the radio interface, including received signal power (RX_power), signal-to-interference-plus-noise ratio (SINR), channel quality indicator (CQI), and scheduling information. Given the inherent non-linearities and dynamic variations in radio environments, traditional analytical models often fall short in accurately mapping these metrics to throughput outcomes. Consequently, supervised machine learning techniques are employed to model the complex statistical relationships, where input vectors comprising historical and current lower-layer measurements are used to estimate instantaneous throughput as the output variable. In this work, we implement Linear Regression (LR) to examine potential linear correlations, Support Vector Regression (SVR) with non-linear kernels such as the Radial Basis Function (RBF) to capture intricate non-linear dependencies, and

Random Forests (RF) to exploit ensemble learning for robust predictions. Furthermore, we introduce lag windows to incorporate temporal dependencies by utilising past metric values, and forecast windows to enable flexible prediction over different time horizons. Model training and hyper-parameter optimisation are conducted using Randomised Search in conjunction with nested K-fold cross-validation to ensure unbiased estimation errors and mitigate overfitting, ultimately enabling precise, generalisable predictions of instantaneous uplink throughput for real-time network optimisation and adaptive resource management.

4. 5G RAN Slicing Enforcement

5G RAN slicing enforcement refers to the set of mechanisms and protocols implemented within the Radio Access Network (RAN) to ensure that the logical network slices-each representing a virtualised, isolated subset of network resources tailored to specific service requirements-are consistently provisioned, maintained, and operated according to their defined Service Level Agreements (SLAs). In the 5G architecture, RAN slicing involves partitioning radio resources, such as spectrum, scheduling capacity, and transmission power, among multiple slices while preserving isolation and quality of service guarantees. Enforcement mechanisms are realised through slice-aware Radio Resource Management (RRM), where functions like admission control, scheduling, and load balancing are augmented with slice-specific policies and prioritisation rules. These include slice-specific configurations for QoS Class Identifiers (QCIs), differentiated treatment of traffic flows, and dynamic resource reservation based on real-time demand fluctuations and network conditions. Furthermore, enforcement relies on standardised interfaces, such as the O-RAN architecture's E2 interface, enabling coordination between the RAN Intelligent Controller (RIC) and distributed units (DUs) for fine-grained control of slice behaviour. Advanced techniques, including machine learning-driven predictive resource allocation and closed-loop optimisation, are increasingly integrated to enhance slice performance and adaptability under diverse traffic scenarios. Through rigorous enforcement of slicing policies, operators can simultaneously deliver heterogeneous services—such as enhanced Mobile Broadband (eMBB), Ultra-Reliable Low Latency Communications (URLLC), and massive Machine-Type Communications (mMTC)-on a shared physical infrastructure while meeting strict performance isolation and SLA compliance requirements.

5. Results

In this study, the feasibility of estimating users' instantaneous uplink throughput in cellular networks based on lower-layer metrics was investigated using a real-time 4G testbed deployed in an anechoic chamber, enabling precise analysis of radio phenomena such as noise, multipath fading, and radio congestion. Measurements were collected at a granularity of 100 ms, constrained by the employed traffic generator. The estimation models were evaluated using three machine learning techniques: Linear Regression (LR), Support Vector Regression (SVR), and Random Forest (RF).

It was observed that radio metrics alone, including SNR, RIP, RSSI, and Rx_power, were insufficient for accurate throughput estimation at very short time scales below 700 ms. However, for forecast windows of 700 ms and longer, acceptable estimation accuracy was achieved using these metrics, particularly in simpler radio environments characterised by linear noise variations. The integration of 43 lower-layer metrics from the eNB significantly improved estimation

performance at small time granularities of 100 ms. When these detailed metrics were utilised, all three machine learning techniques produced accurate throughput estimations under varying radio conditions.

The influence of historical measurements was also examined by introducing a lag window of 100 ms. This led to a slight improvement in estimation accuracy, particularly with the RF model, although the overall benefit was modest and introduced additional complexity without significant performance gains. Among the evaluated techniques, LR consistently delivered estimation accuracy comparable to that of RF and SVR while requiring substantially lower computational time, making it particularly suitable for real-time applications, whereas SVR was consistently the slowest model across all scenarios.

In parallel, experiments were conducted to investigate 5G RAN slicing enforcement through optimisation and heuristic approaches for resource allocation. Although the ESRP models offered optimal resource allocations, their high convergence times limited their suitability for real-time deployment. By contrast, the developed heuristics, namely IMA, HMF, and HSF, achieved competitive performance, particularly in scenarios involving smaller gNB sets, where IMA and HMF attained optimal scores for tied resources and the largest continuous unallocated space (LCUS), with convergence times as low as 10 ms. The heuristics demonstrated robustness, showing little sensitivity to the number of slices served during the allocation process. Additionally, increasing the B size resulted in larger LCUS, enabling more efficient resource allocation strategies and supporting advanced transmission schemes for critical services, albeit sometimes at the expense of total tied resources (TTR).

Overall, the findings confirm that accurate estimation of users' instantaneous uplink throughput at fine time granularities is achievable through the use of detailed lower-layer metrics and suitable machine learning models. Furthermore, the heuristic strategies developed for 5G RAN slicing enforcement demonstrate strong potential for real-time application, ensuring efficient resource allocation while adhering to essential slicing requirements.





Forest (RF) as the underlying machine learning technique.













6. Conclusion and Future Work

This work has addressed the significant challenges facing 4G and 5G cellular networks, which are expected to support diverse sectors requiring sophisticated services and applications with demanding performance criteria and substantial traffic volumes in both uplink and downlink directions. Two principal contributions have been made to meet these challenges.

The first contribution concerns the integration of intelligence into cellular networks through the estimation of users' instantaneous uplink throughput at small time granularities. To this end, a scalable estimation model leveraging machine learning techniques, including Linear Regression, Random Forest, and Support Vector Regression, was developed and evaluated. A real-time 4G testbed was deployed to replicate various radio phenomena, enabling the creation of comprehensive datasets that incorporate cross-layer eNB metrics. The estimation model demonstrated accurate predictions across forecast windows ranging from 100 ms to 1 s, achieving errors below 15% when utilising datasets with extensive lower-layer metrics. It was also concluded that radio measurements alone are inadequate for reliable throughput estimation at small time scales, whereas acceptable estimations can be obtained for larger time granularities from 700 ms onwards.

The second contribution relates to the enforcement of 5G RAN slicing at the resource level from a multi-cell perspective, addressing the stringent requirements associated with RAN slicing, such as orthogonality, satisfaction, scalability, and cooperation enabling. An exact optimisation model using constraint programming was developed to satisfy these requirements, complemented by a 2D bin-packing heuristic to support scalability, albeit at the expense of full cooperation enabling. While the exact model proved effective for large time-scale allocations, it exhibited slow convergence for larger problem instances. Consequently, three heuristics were proposed, prioritising scalability while still addressing all four slicing requirements. Experimental results demonstrated that two of these heuristics delivered strong performance, making them well-suited for real-time RAN slicing deployments.

Looking forward, several avenues for future research and development have been identified. One significant opportunity lies in advancing towards intelligent, proactive systems, where instantaneous throughput estimation could serve as a critical input for sophisticated schedulers and congestion control mechanisms. The growing adoption of network softwarisation and Software Defined Networking (SDN) provides fertile ground for integrating these estimation models into real-time decision-making processes, enhancing the agility and responsiveness of cellular networks.

Another promising direction is the extension of the current experimental framework to a 5G platform, particularly through the development of an open-source testbed based on OpenAirInterface (OAI). Such a platform would facilitate real-time evaluation of slicing heuristics and provide the research community with remote access for validating advanced slicing strategies in versatile 5G scenarios.

Further research is also warranted into refining resource allocation strategies for RAN slicing, especially to balance scalability and cooperation enabling. One potential approach involves designing allocation schemes that, rather than reserving large contiguous unallocated spaces, deliberately create efficiently distributed sparse regions capable of accommodating varying slice demands. This would likely require aggregation heuristics allowing network operators to weight

different requirements according to specific operational priorities. Incorporating traffic forecasting to anticipate slice demand fluctuations could further enhance the effectiveness of such strategies.

Lastly, scaling RAN slicing across broader network domains presents additional challenges and opportunities. Expanding from single SD-RAN domains to large-scale deployments involving multiple SD-RANs necessitates coordinated or cooperative approaches. Coordinated slicing relies on selecting a reference SD-RAN to define resource allocations replicated by neighbouring domains, offering simplicity but potentially limiting adaptability. Alternatively, cooperative slicing involves decentralised decision-making and resource sharing among adjacent SD-RANs, promising greater flexibility and efficiency, albeit at the cost of increased signalling and potential complexity. Developing practical frameworks to support these large-scale cooperative solutions remains a critical area for future exploration.

In summary, while this work has demonstrated the feasibility and effectiveness of both instantaneous uplink throughput estimation and RAN slicing enforcement in cellular networks, further advancements are essential to enable seamless, intelligent, and scalable solutions, particularly as networks evolve towards 5G and beyond.

References

- Chen, S., & Zhao, J. (2020). The requirements, challenges, and technologies for 5G of terrestrial mobile telecommunication. *IEEE Communications Magazine*, 58(5), 36–42. [DOI: 10.1109/MCOM.001.1900417]
- 2. Gupta, A., & Jha, R. K. (2019). A survey of 5G network: Architecture and emerging technologies. *IEEE Access*, 7, 164081–164098. [DOI: 10.1109/ACCESS.2019.2953566]
- Rappaport, T. S., Xing, Y., Kanhere, O., Ju, S., Madanayake, A., Mandal, S., ... Trichopoulos, G. (2019). Wireless communications and applications above 100 GHz: Opportunities and challenges for 6G and beyond. *IEEE Access*, 7, 78729–78757. [DOI: 10.1109/ACCESS.2019.2921522]
- Zhang, J., Chen, S., Han, S., Xu, Y., & Pan, Z. (2021). Research on 5G millimeter-wave technology: Opportunities and challenges. *IEEE Access*, 9, 29523–29544. [DOI: 10.1109/ACCESS.2021.3058827]
- Lin, X., Andrews, J. G., Ghosh, A., & Ratasuk, R. (2020). An overview of 3GPP cellular vehicle-to-everything standards. *IEEE Communications Magazine*, 56(12), 22–28. [DOI: 10.1109/MCOM.001.2000108]
- Ali, A., Qamar, F., Imran, M. A., & Abbasi, Q. H. (2021). A review of 5G spectrum sharing models, architecture, and challenges. *IEEE Access*, 9, 128054–128078. [DOI: 10.1109/ACCESS.2021.3111758]
- Bennis, M., Debbah, M., & Poor, H. V. (2018). Ultrareliable and low-latency wireless communication: Tail, risk, and scale. *Proceedings of the IEEE*, 106(10), 1834–1853. [DOI: 10.1109/JPROC.2018.2867029]
- Taleb, T., Samdanis, K., Mada, B., Flinck, H., Dutta, S., & Sabella, D. (2019). On multiaccess edge computing: A survey of the emerging 5G network edge cloud architecture and orchestration. *IEEE Communications Surveys & Tutorials*, 19(3), 1657–1681. [DOI: 10.1109/COMST.2017.2705720]
- Dai, L., Wang, B., Yuan, Y., Han, S., Chih-Lin, I., & Wang, Z. (2020). Non-orthogonal multiple access for 5G: Solutions, challenges, opportunities, and future research trends. *IEEE Communications Magazine*, 58(7), 86–92. [DOI: 10.1109/MCOM.001.2000022]
- Fouda, M. M., Hamamreh, J. M., & Arslan, H. (2021). Waveform design for 5G and beyond wireless communications: Analysis and comparison. *IEEE Access*, 9, 29335– 29356. [DOI: 10.1109/ACCESS.2021.3058666]
- M. Elsayed, M. Ismail, and Y. Zaki, "A survey of 5G network architecture and emerging technologies: Open issues and future research directions," *Computer Networks*, vol. 194, 2021.
- 12. J. Okwuibe, B. Han, and H. D. Schotten, "Network slicing for 5G: Survey and challenges," *Computer Networks*, vol. 166, 2019.
- 13. M. Alsabah et al., "5G security: A comprehensive review of standards, threats, and solutions," *Computer Standards & Interfaces*, vol. 77, 2021.
- 14. J. Lee and J. Kim, "Artificial intelligence approaches for network optimization in 5G: A survey," *IEEE Access*, vol. 8, pp. 128107–128125, 2020.

- S.-Y. Lien, S.-L. Hung, and Y.-C. Liang, "Massive machine-type communications in 5G and beyond: Recent advances and future trends," *IEEE Internet of Things Journal*, vol. 8, no. 17, pp. 13244–13260, 2021.
- Z. Zhou, H. Zhang, M. Peng, and B. Wang, "Learning-based edge caching in 5G networks: Opportunities and challenges," *IEEE Wireless Communications*, vol. 28, no. 4, pp. 44–50, 2021.
- 17. H. Zhang, J. Chen, B. Zhang, and M. Peng, "Energy-efficient resource allocation for 5G networks: A survey," *IEEE Wireless Communications*, vol. 27, no. 6, pp. 16–23, 2020.
- 18. H. Qureshi, M. A. Khan, and R. Abbas, "Future of mobile networks: 5G architecture, applications, and challenges," *Telecommunication Systems*, vol. 74, pp. 1–17, 2020.
- 19. L. Yang, S. Zhang, and H. Zhang, "Resource allocation in 5G networks: From optimization to learning," *IEEE Wireless Communications*, vol. 28, no. 4, pp. 80–88, 2021.
- 20. M. Khoshnevisan and S. Mirjalili, "A comprehensive review on 5G wireless communication systems: From technologies to applications," *Wireless Personal Communications*, vol. 127, pp. 305–334, 2022.
- M. A. Salehi, H. V. Poor, and W. Saad, "Toward 6G Internet of Things: Recent advances, challenges, and future directions," *IEEE Internet of Things Journal*, vol. 9, no. 2, pp. 1073– 1097, 2022.
- 22. H. X. Nguyen, L. B. Le, and D. I. Kim, "Machine learning for 5G and beyond mobile networks: A survey," *IEEE Access*, vol. 9, pp. 83366–83404, 2021.
- 23. Chih-Lin I and S. Han, "Green 5G networks: Opportunities and challenges," *IEEE Wireless Communications*, vol. 27, no. 2, pp. 16–23, 2020.
- 24. X. Li, H. Zhang, and J. Chen, "5G-enabled industrial IoT: Advances, challenges, and future trends," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 8, pp. 5434–5445, 2021.
- 25. E. Hossain and M. Hasan, "5G cellular: Key enabling technologies and research challenges," *IEEE Instrumentation & Measurement Magazine*, vol. 22, no. 2, pp. 17–25, 2019.
- 26. X. Wang, C. Xu, M. Li, and L. Hanzo, "A survey of 5G non-orthogonal multiple access schemes for massive connectivity," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 4, pp. 2241–2271, 2020.
- 27. R. Hussain and H. Oh, "Autonomous vehicles cybersecurity: Challenges and solutions," *IEEE Communications Magazine*, vol. 57, no. 12, pp. 92–99, 2019.
- 28. S. Ahmadi, 5G NR: Architecture, technology, implementation, and operation of 3GPP New Radio Standards. Academic Press, 2020.
- 29. M. A. Ali, A. Khalid, and I. Hussain, "Evolution of mobile broadband networks from 4G to 5G: A performance comparison," *IEEE Access*, vol. 10, pp. 87622–87636, 2022.
- 30. S. H. Ahmed, H. Gharavi, and B. Chen, "5G-enabled vehicular networks: A review," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 6, pp. 5049–5063, 2021.