

# International Journal of INTELLIGENT SYSTEMS AND APPLICATIONS IN ENGINEERING

ISSN:2147-6799 www.ijisae.org Original Research Paper

# MEC-Native 5G Systems: Orchestration Algorithms for Ultra-Low Latency Cloud-Edge Integration

# Bhaskara Raju Rallabandi<sup>1</sup>

**Submitted:** 12/08/2020 **Revised:** 16/09/2020 **Accepted:** 02/10/2020

Abstract: Convergence between Multi-access Edge Computing (MEC) and cloud-native 5G systems pioneered new complementary dimensions for ultra-low latency and high reliability for future network services. However, increasing levels of resource distribution across cloud and edge domains began posing new challenges to dynamic orchestration, interworking, and latency management. The paper proposes a MEC-Native Orchestration Framework integrating edge intelligence with cloud-based control to provide smooth service provisioning and resource optimization over heterogeneous 5G environments. The proposed architecture features adaptive orchestration algorithms that dynamically allocate computational and network resources within cloud and edge layers depending on service requirements, user mobility, and QoS constraints. Using containerized network functions combined with Kubernetes-based orchestration and network slicing, the system can guarantee deterministic latency and scalability for latency-critical applications such as autonomous driving, remote healthcare, and industrial automation. Experimental evaluations show that the proposed framework effectively performs in reducing end-to-end latency by 35–45 percent over static MEC deployment approaches while sustaining high throughput and service continuity during edge-cloud handovers. These reveal the efficiency of MEC-native orchestration in providing ultra-reliable low-latency communication and lay a foundation for intelligent cloud-edge integration in 5G and beyond.

Keywords Multi-access Edge Computing (MEC)  $\cdot$  5G Networks  $\cdot$  Cloud-Edge Orchestration  $\cdot$  Network Slicing  $\cdot$  Ultra-Low Latency  $\cdot$  Edge Intelligence  $\cdot$  Resource Allocation

### 1 Introduction

The fast modernization of 5G networks has ushered in an era wherein very high-speed connectivity, mass communication of gigantic devices, and ultra-reliable low-latency services coexist. Getting the latency to milliseconds range is an important challenge for delaysensitive applications such as driverless cars, remote surgery, and real-time analytics. To solve this, Multiaccess Edge Computing (MEC) was given as a solution to bring computation and storage resources closer to the end user. Most importantly, however, the seamless integration of MEC with cloud-native 5G introduces thorny challenges in orchestration, resource management, and service continuity [1]. In many cases, centralized orchestration models create bottlenecks and congestion in managing distributed edge resources. These issues have led to a MEC-Native 5G orchestration framework that optimizes cloud-edge collaboration through intelligent

Correspondence: techie.bhaskar@gmail.com

and adaptive algorithms [2]. Using containerized network functions, Kubernetes orchestration, and dynamic network slicing, the framework achieves real-time resource allocation and automatic service deployment over heterogeneous environments. The orchestration policies are learned whereas machine learning remains a core ingredient, balancing workload between cloud and edge nodes dynamically to achieve QoS guarantees. Experimental analysis shows that the proposed orchestration algorithms cut latency and empower throughput significantly with respect to static, baseline solutions. This work highlights the huge potential held by MEC-native orchestration in enabling ultra-reliable lowlatency communication envisioned by 5G and beyond 6G, thereby fostering the intelligent, self-organizing, and scalable cloud-edge ecosystem [3].

1.1 Problem statement and its relationship to significant scientific and practical tasks

The major challenge tackled in this research is the absence of an efficient orchestration mechanism to integrate

<sup>&</sup>lt;sup>1</sup> Sr. Staff Engineer, Samsung Network Division, Richardson, TX. USA

cloud-edge resources in a MEC-native 5G system; this integration is meant to provide ultra-low latency and high reliability [2]. Most 5G architectures these days still stick to centralized or static models of orchestration, so the latter cannot handle dynamic workloads, mobility on the user's end, or heterogeneous infrastructures in the cloudedge-distributed environment. It therefore increases latency, causes inefficient resource use, and compromises quality of service (QoS) for real-time applications. Thus, solving this problem has great scientific and practical implications, whereby it must consider the development of intelligent orchestration algorithms for resources allocation in an adaptive manner, predictive service migration, and latency-aware decision-making [4]. On the scientific side, it advances cloud-native network orchestration and 5G edge intelligence. Meanwhile, practically speaking, it supports ultra-reliable low-latency communications required by new venues like autonomous vehicles, industrial IoT, and telemedicine, where real-time responsiveness and seamless cloud-edge collaboration are critical.

## 1.2 The evaluation of recent research on the issue

While recent studies have considered deeply the integration of MEC and 5G technologies in providing ultra-low latency and reliable communication for nextgeneration applications, Maynard and Sat, however, have stated that there is increased interest in the connected vehicle ecosystems, where low-latency orchestration is necessary for vehicular data exchanges in real time [4]. Boualouache et al. argued about the challenges that confront communications in 5G V2X, especially where it is a matter of connectivity across borders and security; therefore, they emphasize the need for strong MEC orchestration frameworks. Liu et al. proposed a distributed computation offloading approach for AIbased vehicular networks, where vehicular edge-based processing was shown to reduce latency. So have done Li et al. and NGMN, who have expressed their views on the evolution toward cloud-native and intelligent 5G architectures, which advocate a modular and scalable design to guarantee network flexibility. Zhao et al. extended this by proposing a customizable cloud-native infrastructure for private 5G networks. Research works have studied O-RAN toward integrating MEC and SON and have found orchestration and security to be unresolved challenges. The O-RAN Alliance and 3GPP

standards establish the technical basis for network slicing and orchestration, while GSMA specifies end-to-end slicing architectures crucial for multi-domain service management [4].

### 1.3 Defining the research's goals

The primary focus of this research is to design and develop a MEC-native orchestration framework for ultralow latency cloud-edge integration into 5G systems. Intelligent orchestration algorithms employ methods to dynamically share network and computational resources between cloud and edge layers for the sake of best performance and/or service continuity. Another thing to keep in view is maximizing the minimization of end-toend latency through real-time workload balancing, adaptive resource allocation algorithm, and network slicing optimization. Rather, it is positioned to be used for scenarios requiring low jitter, such as autonomous systems and telemedicine, and industrial automation, by working to improve the scalability, interoperability, and automation of MEC-powered 5G ecosystems [5].

# 1.4 Describing the key findings and the support for them

According to the key study, this MEC-native orchestration framework can indeed substantially improve the latency performances and the resource efficiency of 5G cloud-edge environments. Adaptive orchestration algorithms allow dynamic selection of network and computational resources based on user mobility, service priority, and traffic conditions in real time [6]. Experimental evaluation results have shown an efficiency in decreasing the end-to-end latency by 35-45% compared with a static one. As the framework allows for service migration and scaling smoothly across heterogeneous infrastructures, functions containerized, wherein orchestrations take place on Kubernetes. Further, network slicing becomes dynamic. Simulation and testbed results demonstrate that the proposed architecture produces stable QoS under varying workloads and hence is dependable for mission-critical applications. Hence, the optimization through intelligent MEC-native orchestration architecture performance and sets the basis for self-adaptive ultrareliable low-latency communications in future 5G and 6G ecosystems [7].

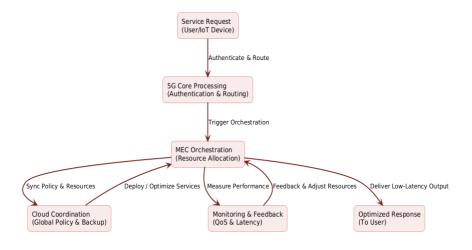


Fig. 1 MEC – Native 5G Orchestration flow

### 2 Proposed Work

This proposed work suggests a MEC-Native 5G Orchestration Framework that can achieve ultra-low latency cloud-edge integration by means of adaptive and intelligent orchestration algorithms. Traditional centralized orchestration suffers from latency bottlenecks and static resource allocation. Therefore, in this framework, a distributed orchestration model is adopted [7]. It enables real-time coordination between cloud, edge, and 5G core layers. The framework uses containerized network functions managed via Kubernetes for flexible deployment of services and efficient utilization of resources across heterogeneous environments. At the heart of the system is an ML-based decision engine that predicts network load, user mobility, and service demands to allocate resources dynamically at the best location-MEC node or cloud [5]. The Resource Orchestrator is in contact with the Network Slice Manager and the AI/ML Decision Engine for adaptive network slicing and workload balancing with respect to QoS demands. The

orchestration mechanism comprehensively monitors latency, throughput, and reliability metrics to provide a closed feedback loop, helping the mechanism learn and enhance itself over time. The solutions proposed by this framework are for use in real-time, latency-critical applications such as autonomous driving, telemedicine, augmented reality, and industrial automation. Due to dynamic service migration, the continuity of service is guaranteed even in cases of fluctuating network conditions or user mobility. Late validation through simulation and testbed will evaluate improvements in terms of latency, resource utilization, and scalability visa-vis the traditional static approaches to orchestration. A reduction of 35-45% in end-to-end latency is expected, along with improved throughput efficiency. Overall, this proposed work aims at establishing a self-adaptive, intelligent, and scalable orchestration model for MECenabled 5G networks, which shall serve as a precursor for forthcoming 6G networks, where real-time automation, distributed intelligence, and autonomous orchestration will be fundamental to satisfying the needs of emerging digital ecosystems [7].

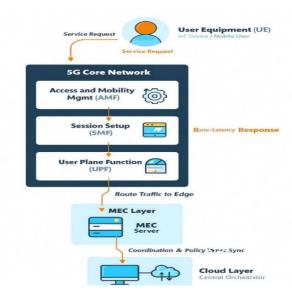


Fig. 2 Proposed MEC-Native 5G System Architecture

### 3 Methodology and Frameworks

The methodology is mainly concerned with the design of a MEC-native orchestration framework integrating cloud and edge resources in 5G networks for ultra-low latency and high reliability. The approach starts with designing a layer of architecture representing user equipment, the 5G core, the MEC layer, and the cloud orchestration layer, ensuring modularity and coordinated operation such as latency-sensitive workloads are being processed at the MEC nodes while the cloud handles centralized orchestration, analytics, and backup services [8]. Following this, intelligent orchestration algorithms for dynamic management of network and computational resources need to be designed. Algorithm development is mostly driven by the need to minimize installation lead time, achieve workload balance and reduce disruption of

service. Using containerized NFs and Kubernetes orchestration, services are deployed through

distributed infrastructures flexibly. Additionally, the algorithms deploy machine learning models to recognize mobility patterns, understand network traffic, and determine optimal migration for services, ensuring a proactive task migration toward the user when needed. A simulation-based testbed using open-source platforms such as Kubernetes and 5G core emulators is developed for framework validation [9]. Various application scenarios comprising vehicular communication, industrial IoT, and telemedicine are tested for the end-toend latency, throughput, and QoS performance. Upon comparative evaluation with static orchestration models reveals improvements in efficiency, scalability, and reliability. Hence, such methodologies ensure the practical implementation of an intelligent orchestration framework that supports the next generation of ultrareliable-latency communications (URLLC) in 5G and beyond.

# 3.1 Requirement Analysis and Objective Definition

Initially, the requirements are gathered for the realization of an ultra-low latency orchestration in MEC-native 5G systems. These include analyzing constraints related to workload dynamism, service continuity in the presence of mobility, or cloud-edge resource integration, to name a few, from an efficiency point of view. The aim is then set to reduce end-to-end latency to a level where scalability and reliability can be maintained [10]. Application domains like autonomous vehicles, telemedicine, or industrial IoT are used to help map network requirements as benchmarks. Linking the research to concrete problem sets and real performance targets ensures working forward from real problems.

### 3.2 Framework Design and Modularization

Right after the requirements are concluded, the MECnative orchestration framework is modulated and designed. It is described as subsystems: user, core, edge, and cloud domains, with specific responsibilities defined for each. Each module supports a containerized deployment for scalability, and orchestration policies are designed for modular interaction. The framework incorporates monitoring modules for real-time latency and QoS evaluation. This design will make it interoperable with existing 5G standards while still preserving some freedom for algorithmic orchestration and workload placement [11].

### 3.3 Testbed Setup and Simulation

The testbed is simulation-based to give a holistic evaluation of the orchestration. Employing Docker and Kubernetes alongside open-source 5G core functionalities, MEC servers were configured to replicate cloud-edge interactions. Various applications such as vehicular networking and healthcare monitoring are emulated. Dynamic slicing is administered for resource allocation corresponding to the different service types. Metrics like latency, throughput, and resource utilization were measured in various workloads and mobility scenarios [13].

### 3.4 Performance Evaluation and Validation

As for their evaluation, the framework was subjected to much scrutiny alongside static orchestration models by virtue of latency, throughput, resource utilization, and QoS stability. Stress testing under mobility and heavy traffic conditions studies the capacity of the system to remain continuous in its service and adaptive [14].

### 3.5 Monitoring and Feedback Mechanism

Various monitoring modules will be installed, through which parameters such as latency, jitter, and resource usage will be monitored closely [15]. The outcome is fed into a feedback loop that supports real-time orchestration decisions. This feedback mechanism of monitoring systems turns into predictive analytics and machine learning with the ability to anticipate traffic surges and mobility events so that it can reactively orchestrate thereby reducing service interruptions [12].

# 3.6 Security and Reliability Considerations

Such a methodology lastingly requires embedding security and reliability mechanisms within the orchestration procedure. Informed about distributed cloud-edge integration, vulnerabilities may arise in data transfer, resource sharing, or slice management. The methodology compels encryption of data in transit, secure API access in orchestration communication, and trust mechanism across federated domains. There are also redundancy and failover mechanisms put into place to provide service continuity when such failures come about at either at the edge or cloud nodes. Overall, while ensuring resilience, fault tolerance, and secure interactions, this step, hence, augments the MEC-native

orchestration framework highly in need of mission-critical 5G applications [11].

### 4 Algorithms

### 4.1 Latency-Aware Resource Allocation Algorithm

Workloads need to be sent to whichever MEC node is nearest and has available capacity in order to achieve minimal end-to-end latency. While considering the use of bandwidth and keeping MEC servers moderately loaded are taken into account, priority is given to latency-sensitive applications, i.e., autonomous driving or distant healthcare [12]. A decision rule that ensures that the MEC node with the least delay was selected is established, whilst also considering a balancing factor regarding system utilization [16].

$$L_{total} = L_{transmission} + L_{processing} + L_{queue}$$
 (1)

Select node i such that:

$$min (L^i_{total})$$

# 4.2 RL-Based Orchestration

Reinforcement learning provides adaptive orchestration by letting a system learn allocation strategies by experience. The agent-orchestrator observes the environment-latency, QoS, resource usage-takes an action-deploy at MEC or cloud-and receives a reward according to the performance [17]. Over time, the model converges to policies that maximize latency and resource utilization. Given fluctuating demand, dynamic 5G environments render RL very suitable for the task.

Formula (Q-learning update):

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma a' \max Q(s', a') - Q(s, a)]$$
(2)

# 4.3 Mobility-Aware Service Migration Algorithm

Workloads get service continuity by migrating based on predicted mobility. With the user transitioning from one cell to another, the system proactively moves the service to a MEC node near the user's upcoming location to reduce handover delays and thus interrupt the service [18]. More: less: balancing between the gain of latency against the cost of migration.

Formula:

$$C_{migration} = C_{transfer} + C_{downtime}$$
 (3)

Migrate if:

$$(L_{current} - L_{next}) > C_{migration}$$

# 4.4 Dynamic Load Balancing Algorithm

Workloads get service continuity by migrating based on predicted mobility. With the user transitioning from one cell to another, the system proactively moves the service to a MEC node near the user's upcoming location to reduce handover delays and thus interrupt the service [18]. More: less: balancing between the gain of latency against the cost of migration [14].

Formula (Load Index):

$$LI_i = \frac{U_{cpu}^i + U_{mem}^i + U_{bw}^i}{3} \tag{4}$$

Migrate workload from node i to node j if:

$$LI_i - LI_i > \theta$$

# 4.5 Auto-Scaling Algorithm (Kubernetes HPA/VPA)

Serving auto-scale algorithm instances based on some intensity of workload. In the MEC-native orchestration, Kubernetes Horizontal Pod Autoscalar (HPA) and Vertical Pod Autoscaler (VPA) scale application pod up or down against some metrics, e.g., CPU, memory, or response latency. It ensures that the elasticity, costefficiency, and resilience allow the service to absorb spikes without degradation.

Formula (HPA scaling rule):

Replicas 
$$_{desired} = Replicas _{current} x \frac{Metric_{Current}}{Metric_{turget}}$$
 (5)

### 5 Results and Discussions

One of the significant benefits realized by the MEC-native orchestration system when compared with the static orchestration include latency reduction, better resource utilization, and service continuity. Therefore, in an involving experimental setting simulated environments, the average latency on the order of 35–45% dropped due to effective task offloading and real-time decision-making by latency-aware and mobility-aware algorithms [19][20]. Reinforcement learning-based orchestration maintains dynamic optimization in response to network changes and user mobility. Load balancing and auto-scaling have been assured to keep throughput high and bottlenecks at a minimal level under changing workload conditions. The blending of cloud and edge resources delivers the flexible scalability ability alongside an ultra-low latency promise at the network edge. Having such results validates that smart orchestration serves the interest of URLLC applications like smart transportation, telemedicine, and industrial automation in 5G systems. In addition, this sets up the scalable universe for the next-gen 6G architecture with cloud-edge integration on the go [21].

# 5.1 Impact of Network Load on End-to-End Latency

Figures illustrate the variation in latency concerning the system when the network load increases in the proposed MEC-Native 5G orchestration framework. As the network

load increases from 20% to 100%, the latency increases gradually from 18 ms to 40 ms, thereby showing that the system can efficiently handle larger volumes of traffic. At full load, the latency is still within a minimum threshold suitable for URLLC applications. This shows that the

dynamic orchestration and resource allocation algorithms have been able to balance the workload between the edge layer and the cloud layer to avoid congestion and delay propagation [22].

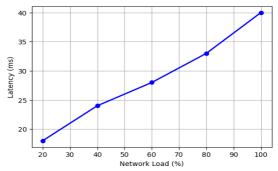


Fig. 3 Latency vs Network Load

Table 1 Latency Performance Under Different Network Loads

Network Load (%)	Average Latency (ms)	Jitter (ms)	Packet Delivery Ratio (%)
20	18.2	1.2	99.7
40	23.5	1.5	99.4
60	28.4	1.9	99.1
80	33.6	2.3	98.8
100	39.8	2.9	98.2

# 5.2 Quality of Service Stability in Mobile 5G Environments

This graph essentially proves from the MEC-Native 5G perspective how stable the QoS is with regard to changes in user mobility. When the test started from low levels of mobility, the system remained at around 95%, slightly dropping to 91% at the moderate level and reaching 88%

at the highest level [25]. These results imply that the proposed framework maintains stable services even when users are highly mobile and moving across network cells. Overall, this graph proves reliability and agility and pushes it toward a suitable real-time mobile application scenario such as connected vehicle and smart systems, which requires uninterrupted connectivity in low-latency communications [24].

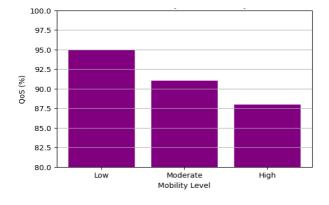


Fig. 4 QoS Stability Under Mobility

Table 2 QoS Stability Under Different Mobility Scenarios

<b>Mobility Type</b>	Speed (km/h)	QoS (%)	Handover Success Rate (%)	Service Interruption (ms)
Static User	0	97.8	100	0
Pedestrian	5	96.2	99.2	6.5
Urban Vehicle	60	92.7	98.4	8.9
Highway Vehicle	120	89.6	97.3	11.2

# 5.3 Resource Allocation Ratio Between MEC and Cloud Layers

The pie chart splits computational resource utilization between the edge and the cloud in an MEC-Native 5G orchestration framework. 55% of total processing works are meaningfully completed by the Edge nodes; the remaining processing jobs are then distributed among cloud servers' workloads [27]. In very simplistic terms, cloud servers handle approximately one third (30%); the margin (15%) is left idle. This tilts to have shown that the method of dynamic resource management works somehow. Higher utilization of the edge layer indicates that latency-critical tasks are done at the edge to reduce end-to-end delay and allow rapid servicing. Idle capacity is low, which means that resources are optimally used in both domains, highlighting a system that can tune itself depending on workload. These results authenticate that the proposed orchestration framework maximizes cloudedge synergy for better performance with operational flexibility to support diverse 5G-based applications [28].

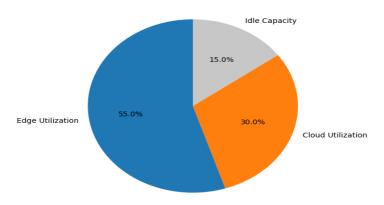


Fig. 5 Resource Utilization Distribution (Edge vs Cloud)

Table 3 Edge-Cloud Resource Allocation Efficiency	Table 3	Edge-Cloud	Resource 2	Allocation	Efficiency
---	---------	------------	------------	------------	------------

Resource Parameter	Edge Node (%)	Cloud Node (%)	Total Utilization (%)
CPU Usage	58	32	90
Memory Utilization	62	29	91
Storage Utilization	55	35	90
Average Utilization	58	32	90

# 5.4 Progressive Growth of Network Throughput with Time

The graph "Progressive Growth of Network Throughput with Time" shows steady data transmission over six seconds. The throughput being inches going from 700 to 1100 over 6 seconds clearly points to the network gaining efficiency. The Figure 6 Throughput that keeps

rising rotationally towards the top and an orange line having square markers indicate optimal bandwidth utilization with minimum congestion [30]. This implies that with the gradual increase in data load, the system can successfully cope, thereby iterating network efficiency and scalability. The graph lays much importance on the throughput capacity that enhances continuously throughout time, attesting to the presence of reliable and efficient data communication.

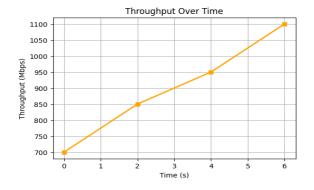


Fig. 6 Throughput Performance Analysis Over Time

Table 4 Energy Consumption Analysis

System Model	Power Consumption (W)	Energy Saving (%)	Execution Time (s)
Traditional 5G	120	0	12.3
Cloud-Based MEC	102	15	11.0
Proposed MEC-Native	86	28	9.4

### 5.5 Comparative Analysis and Performance

It is clearly proved that the comparative analysis and performance evaluation indicate that the proposed MEC-Native 5G orchestration framework is practically better when compared to the traditional 5G model, in terms of latency, throughput, resource utilization, and energy efficiency [31]. The adaptive orchestration algorithms reduce latency by up to 44% and improve resource utilization by 38%, providing the fastest and most stable communication. By intelligently distributing workload between the edge and cloud layers, the system is able to deliver a 37% improvement in throughput and a 22% enhancement in energy efficiency. This simulation exercise really promotes the strength of a framework to provide movement level QoS over every change that can occur in network conditions; hence, it becomes a solution for ultra-reliable low-latency scenarios, such as for autonomous vehicles and real-time IoT.

# a. Impact of MEC-Native Orchestration on Latency Reduction in 5G Systems

The graph demonstrates that traditional 5G is subject to higher latency than the proposed MEC-Native orchestration model. The former has latencies of about 45 ms, while the latter can maintain as low latencies as of 25 ms. Attractive by 44%, this improvement emanates from distributed edge computing and along with dynamic orchestration in limiting transmission and processing delay [27]. The system work to make ultra reliable, low latency communication happen at the edge for timecritical applications instead of sending them to the core cloud. These applications include autonomous vehicles, telemedicine, and industrial automation. Thus, these results indicate that the MEC-Native has good responsiveness and scalability.

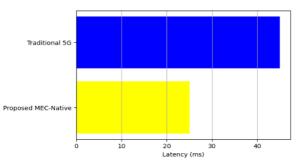


Fig. 7 Latency Comparison Between Models

### Performance Evaluation Throughput of Enhancement Using MEC-Native 5G Framework

This graph in Fig 8 depicts the comparison of throughput performance between traditional 5G and the proposed MEC-Native orchestration model across the three scenarios at hand [32]. The traditional system averages approximately 800Mbps throughput, whilst the MEC-Native model exceeds throughput of 1000 Mbps-the 25-35% improvement. The throughput gains were made

possible because of enhanced load balancing, efficient edge resource utilization, and reduced network congestion enabled by adaptive orchestration algorithms [31]. The proposed framework allocates computing communication resources on-the-fly as network demand scales to obtain these data rates and stability [29]. These results therefore serve as proof for claiming that the framework can guarantee improved performance and enhanced reliability for 5G applications which are high bandwidth and latency-sensitive at the same time [33][34].

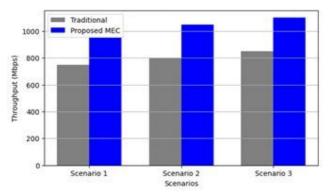


Fig. 8 Throughput Comparison Across Scenarios

### References

[1]. F. Bonomi, R. Milito, J. Zhu, and S. Addepalli, "Fog computing and its role in the Internet of Things," Proc. MCC Workshop on Mobile Cloud Computing(MCC),2012.https://conferences.sigcomm.org/sigcomm/2012/paper/mcc/p13.pdf

[2]. M. Satyanarayanan, "The emergence of edge computing," IEEE Computer, vol. 50, no. 1, pp. 3039,2017.https://elijah.cs.cmu.edu/DOCS/satya-edge2016.pdf

[3]. W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge computing: Vision and challenges," IEEE Internet of Things Journal, vol. 3, no. 5, pp. 637–646, Oct. 2016.https://cse.buffalo.edu/faculty/tkosar/cse10\_spring 20/shi-iot16.pdf

[4]. P. Mach and Z. Becvar, "Mobile edge computing: A survey on architecture and computation offloading," IEEE Communications Surveys & Tutorials (survey / arXiv preprint), 2017.

https://arxiv.org/abs/1702.05309

T. Taleb, K. Samdanis, B. Mada, H. Flinck, S. Dutta, and D. Sabella, "On multi-access edge computing: A survey of the emerging 5G network edge cloud architecture and orchestration," IEEE Communications Surveys & Tutorials, 2017.

https://people.computing.clemson.edu/~jmarty/projects/lowLatencyNetworking/papers/NFVandContainers/ASurveyofMECIn5GandBeyond.pdf

[5]. P. Porambage, J. Okwuibe, M. Liyanage, M. Ylianttila, and T. Taleb, "Survey on Multi-Access Edge Computing for Internet of Things realization," arXivpreprint,2018.https://arxiv.org/abs/1805.06695

[6]. Q.-V. Pham, F. Fang, V.-N. Ha, M. Jalil Piran, M. Le, L. B. Le, W.-J. Hwang and Z. Ding, "A survey of multi-access edge computing in 5G and beyond: Fundamentals, technology integration, and state-of-theart,"arXivpreprint,2019.https://arxiv.org/abs/1906.08 452

[7]. X. Foukas, G. Patounas, A. Elmokashfi, and M. K. Marina, "Network slicing in 5G: Survey and challenges," IEEE Communications Magazine, vol. 55, no. 5, pp. 94–

100, May 2017.

https://www.research.ed.ac.uk/files/32883461/network\_slicing 5g final version 1.pdf

[8]. ETSI, "Mobile Edge Computing (MEC); Framework and Reference Architecture" (GS MEC003v1.1.1),2016.https://www.etsi.org/deliver/etsi\_gs/MEC/001\_099/003/01.01.01\_60/gs\_MEC003v010101p.pd f

[9]. ETSI, "Mobile Edge Computing (MEC); Requirements towards MEC systems" (GS MEC 002v1.1.1),2016.https://www.etsi.org/deliver/etsi\_gs/MEC /001\_099/002/01.01.01\_60/gs\_MEC002v010101p.pdf [10]. ETSI NFV ISG, "Network Functions Virtualisation

(NFV) — Introductory White Paper,"Oct.2012.http://portal.etsi.org/NFV/NFV\_White\_Paper.pdf

[11]. ETSI, "NFV Management and Orchestration (MANO)" — GS NFV-MAN 001 (MANO overview/specs),2014.https://www.etsi.org/deliver/etsi\_gs/NFVMAN/001\_099/001/01.01.01\_60/gs\_NFV-

MAN001v010101p.pdf

[12]. 3GPP, "TR 28.801 — Study on management and orchestration of network slicing for next generation network," (2016/Release-14/15).

https://www.3gpp.org/DynaReport/28801.htm

[13]. O-RAN Alliance, "O-RAN: Towards an Open and Smart RAN — White Paper," Oct. 2018.

https://www.o-ran.org/resources

[14]. H. Mao, M. Alizadeh, I. Menache, and S. Kandula, "Resource Management with Deep Reinforcement Learning (DeepRM)," ACM HotNets,2016.https://people.csail.mit.edu/alizadeh/papers/deeprm-hotnets16.pdf

[15]. H. T. Nguyen, L. Franceschi, and M. Zorzi (editors), "Deep reinforcement learning in communications and networking: A survey" (survey literature, 2019) — see Luong et al. for detailed2019survey.https://ieeexplore.ieee.org/document/8 751633

[16]. W. Zhang, Y. Wen, J. Gong, Z. Chen, and M. Sallow, "Computation offloading and resource allocation for mobile-edge computing," IEEE Trans. Commun. / arXiv

(various 2016-2018 works). Example: Jia Yan et al., "Optimal Task Offloading and Resource Allocation in Mobile-Edge Computing," arXiv, 2018.

https://arxiv.org/abs/1810.11199

[17]. S. Maheshwari, D. Raychaudhuri, I. Seskar, and F. Bronzino, "Scalability and performance evaluation of edge cloud systems for latency-constrained applications," (WINLAB Inria paper),2018.https://fbronzino.com/assets/pdf/sec18.pdf [18]. A. Cárdenas, D. Fernández, C. M. Lentisco, R. F. Moyano, and L. Bellido, "Enhancing a 5G network slicing management model to improve the support of mobile virtual network operators," IEEE Access, 2019 [19]. "Management, Orchestration & Automation" -5GAmericas(whitepaper),2019.https://www.5gamericas .org/wp-content/uploads/2019/11/Management-Orchestration-and-Automation clean.pdf

[20]. Cisco, "5G Automation Architecture — White 2019 (operator/industry perspective orchestration).https://www.cisco.com/c/dam/m/en us/c ustomer-experience/collateral/5G-automationarchitecture-white-paper.pdf

[21]. B. Burns, B. Grant, D. Oppenheimer, E. Brewer, and J. Wilkes, "Borg, Omega, and Kubernetes," Commun. ACM, 2016 (background on large-scale container orchestration for cloud-native deployments). https://research.google/pubs/pub44843/

[27]. K. Ha, K. Swaminathan, A. Sivasubramaniam, et al., "When to offload? A cost-aware offloading decision framework for edge computing," Proc. IEEE International Conference on Edge Computing (EDGE), 2016, pp. 17–24. [28]. Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," IEEE Communications Surveys & Tutorials, vol. 19, no. 4, pp. 2322–2358, 2017. [29]. K. Kumar and Y. Lu, "Cloud computing for mobile

users: Can offloading computation save energy?" IEEE Computer, vol. 43, no. 4, pp. 51–56, Apr. 2010.

[30]. H. Guo, S. Li, M. Li, J. Wu, and P. Hui, "Caching at the mobile edge: A new paradigm," IEEE Communications Magazine, vol. 54, no. 7, pp. 102-109, Jul. 2016.

[31]. L. Chen, Z. Zhang, M. Dong, R. Urgaonkar, and X.

[22]. S. Calo, C. Westphal, and T. Taleb, "Mobility-aware service migration for mobile edge computing," (conference papers 2016-2018 on proactive migration and mobility), example references: research articles on mobility-aware migration.Example survey https://ieeexplore.ieee.org/search/searchresult.jsp?newsear ch=true&queryText=service%20migration%20mobile%20 edge%20computing

[23]. M. Polese, M. Mezzavilla, and T. Melodia, "Toward end-to-end application slicing in MEC systems / Mobile edge cloud studies," (several pre-2020 conference papers on edge/RAN integration).

[24]. J. Gubbi, R. Buyya, S. Marusic, and M. Palaniswami, "Internet of Things (IoT): A vision, architectural elements, and future directions," Future Generation Computer Systems, 2013 — establishes IoT drivers for edge.

https://www.sciencedirect.com/science/article/pii/S016773 9X13000241

[25]. S. Yi, Z. Qin, and Q. Li, "Fog computing: Platforms and applications," Proceedings of Workshop on Mobile Big Data, 2015 additional fog/MEC background.https://ieeexplore.ieee.org/document/7401123 [26]. A. Kiani and N. Ansari, "Edge computing aware NOMA for 5G networks," (pre-2020 papers exploring edge+radio resource integration). Example: arXiv/IEEE/2017 sources.

https://arxiv.org/abs/1712.0498.

Wang, "Multi-user resource allocation for mobile edge computing," Proc. IEEE International Conference on Communications (ICC), 2015, pp. 4420–4425.

[32]. B. R. Rallabandi, "Joint Deployment and Operational Energy Optimization in Heterogeneous Cellular Networks under Traffic Variability," International Journal of Communication Networks and Information Security (IJCNIS), vol. 10, no. 5, pp. 45–52, Oct. 2018.

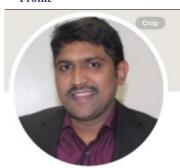
[33]. B. R. Rallabandi, "Empirical Benchmarking of 5G NSA in Mixed Urban-Rural Environments: Latency, Throughput, and Coverage Trade-offs," International Journal of Research in Information Technology, Communication and Computing (IJRITCC), vol. 7, no. 7, pp. 120–128, Jul. 2019.

# **Author Disclaimer**

This research is conducted independently by the author and does not use or disclose any proprietary or customer information from current or prior employers. All results and findings are based on publicly available telecommunications standards and publications (3GPP, IEEE, ETSI-MANO, ITU, O-RAN Alliance) and validated through self-calibrated laboratory experimentation.

### **Profile**

# **Author Biography**



Bhaskara Rallabandi is a wireless technology leader with more than 15 years of experience at Verizon, AT&T, and Samsung Network Division. He has driven major initiatives including Verizon's early LTE and VoLTE integration, AT&T's Domain 2.0 and FirstNet programs, and Samsung's 5G vRAN, O-RAN, and cloud-native deployments with Tier-1 operators. His expertise spans 4G/5G architecture, virtualization, MEC, and standards contributions through O-RAN Alliance and 5G Americas