5G MEC Computation Handoff for Mobile Augmented Reality

Pengyuan Zhou University of Helsinki pengyuan.zhou@helsinki.fi

Sasu Tarkoma University of Helsinki sasu.tarkoma@helsinki.fi Benjamin Finley University of Helsinki benjamin.finley@helsinki.fi

Jussi Kangasharju University of Helsinki jussi.kangasharju@helsinki.fi

Pan Hui University of Helsinki Hong Kong University of Science and Technology pan.hui@helsinki.fi Xuebing Li Aalto University xuebing.li@aalto.fi

Mostafa Ammar Georgia Institute of Technology ammar@cc.gatech.edu

ABSTRACT

The combination of 5G and Multi-access Edge Computing (MEC) can significantly reduce application delay by lowering transmission delay and bringing computational capabilities closer to the end user. Therefore, 5G MEC could enable excellent user experience in applications like Mobile Augmented Reality (MAR), which are computation-intensive, and delay and jitter-sensitive. However, existing 5G handoff algorithms often do not consider the computational load of MEC servers, are too complex for real-time execution, or do not integrate easily with the standard protocol stack. Thus they can impair the performance of 5G MEC.

To address this gap, we propose *Comp-HO*, a handoff algorithm that finds a local solution to the joint problem of optimizing signal strength and computational load. Additionally, *Comp-HO* can easily be integrated into current LTE and 5G base stations thanks to its simplicity and standard-friendly deployability. Specifically, we evaluate *Comp-HO* through a custom NS-3 simulator which we calibrate via MAR prototype measurements from a real-world 5G testbed. We simulate both *Comp-HO* and several classic handoff algorithms. The results show that, even without a global optimum, the proposed algorithm still significantly reduces the number of large delays, caused by congestion at MECs, at the expense of a small increase in transmission delay.

1 INTRODUCTION

Cellular networks are a vital part of modern society with novel network technologies enabling an expanding array of use cases from simple NB-IoT sensors to immersive mobile virtual reality. Fifth-generation mobile networks (5G) specifically provide support for much higher frequencies (up to 52.6 Ghz) with larger bandwidths (up to 400 Mhz) and lower radio access network delay (around 10 ms) in comparison to LTE.

Relatedly, Multi-access Edge Computing (MEC), another novel networking technology, supports deploying compute nodes near existing network nodes in the mobile network structure (often as servers co-located with base stations in the radio access network). Thus user applications that require computation can lower total delay by sending the computation request to a physically and hopwise closer compute node rather than to a remote cloud server [34,

45]. MEC is particularly suitable for computation-intensive and delay and jitter -sensitive applications such as mobile augmented reality (MAR) [11, 21, 43].

MEC has recently been standardized by ETSI thus detailing the potential for the tight integration of MEC and 5G technologies. In such a context, the handoff process [37] between base stations is an important and growing concern. This is because such a handoff often means the user will also be served by a different MEC server. Under the assumption that different MEC servers have varying capabilities and loads, the user application performance will thus be affected by the capability and load of the new MEC server. However, current handoff decisions are typically based primarily on communication measurements such as signal strength, without concern for the status of MEC servers. Therefore, a handoff that improves signal strength could still reduce overall application quality due to load differences of MEC servers. This concern is especially important in 5G networks given their typically small cell sizes (<500m) which implies more frequent handoffs. We denote this issue as the MEC HO problem.

To address this issue, in this work we propose *Comp-HO*, a low-complexity, stand-friendly handoff algorithm jointly considering received signal quality from base stations and the computational loads of the co-located MEC servers. In other words, when the signal strength degrades sufficiently or the serving MEC server is sufficiently overloaded, the base station initiates a handoff and re-assigns the communication and computation processes to another base station and MEC server. We focus specifically on the MAR use case, which has strict real-time quality requirements. Furthermore, we show that our approach strikes a good balance between the signal strength and computational load concerns with large numbers of MAR users.

Specifically, we first develop a MAR prototype and deploy the prototype in a 5G MEC testbed (shown in Figure 1) to measure baseline performance with a traditional handoff algorithm. We then utilize the measurement results to drive a custom MEC-enabled NS-3 simulator and compare the *Comp-HO* algorithm to classic handoff algorithms on larger scale network simulations.

 $^{^1\}mathrm{Unfortunately},$ the commercial 5G base station does not allow reprogramming so we cannot test the Comp-HO algorithm in the testbed.

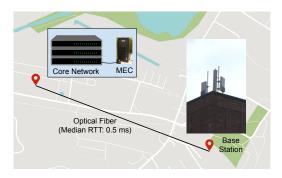


Figure 1: The 5G testbed is composed of the core network and the base station, which are interconnected via an optical fiber. The core network hosts the MEC via its NG6 interface, which is designed for connecting to the data network. The base station has two 5G antennas and one LTE antenna (middle).

To the best of our knowledge, ours is one of the first efforts to address the 5G *MEC HO* problem. Our contributions are threefold:

- (1) Defining the MEC HO problem in 5G and proposing an algorithmic solution, Comp-HO, that jointly considers signal strength and computational load. Comp-HO is simple and standard-friendly and outperforms traditional algorithms with only a small transmission overhead.
- (2) Measuring baseline performance with a MAR prototype in a real-world 5G MEC testbed with a traditional handoff algorithm. This baseline helps guide the parameter selection for the custom NS-3 simulations.
- (3) Carrying out reproducible and measurement-based simulations to evaluate *Comp-HO* at scale using a custom MEC-enabled NS-3 network simulator [4].

The remainder of this paper is organized as follows. Section 2 presents related work in the areas of edge offloading and handoff algorithms. We describe mathematically the *MEC HO* problem in Section 3 and propose *Comp-HO* algorithm in Section 4. Section 5 details the baseline measurements with the 5G testbed. Section 6 presents the simulation setup and results. We discuss the potential future directions of the work in Section 7 and conclude in Section 8.

2 RELATED WORK

Handoffs can be classified into horizontal and vertical [32, 35], or hard and soft [10, 36]. In this work, we focus on horizontal hard handoff since we target a pure 5G networking environment and hard handoff is more common with LTE [16].

Traditional works largely make handoff decisions based on UE measurements consisting of signal quality indicators such as received signal strength [26, 39], signal-to-interference-plus-noise ratio (SINR) [12, 40] and reference signal received quality (RSRQ) [19, 41]. However, given the novelty of MEC, the mentioned algorithms only focus on transmission metrics without taking edge computation into consideration. The lack of such algorithms could be partly because the radio access network delay of LTE does not support some MEC-assisted applications [15] compared to novel standards like 5G, therefore lessening the motivation for LTE + MEC research.

As MEC and 5G techniques evolve, there has been some MEC-aware handoff research. Nasrin et al. [27] propose a handoff algorithm that jointly considers signal quality and computational loads. Sardellitti et al. [33] and Mao et al. [23] focus on the joint optimization of radio and computational resources considering energy consumption and user experienced delay.² Basic et al. [5] propose a fuzzy logic handoff algorithm that selects a target node based on bandwidth, processor, and delay parameters of edge servers. Emara et al. [13] and Li et al. [20] both propose to improve handoff algorithms by considering MEC load in 5G heterogeneous networks and cloud radio access networks, respectively. Finally, Zhang et al. [46] detail a UE-based offloading algorithm for dense networks with MECs that considers transmission delay, processing delay, and an energy constraint under information uncertainty. Though in contrast to our work, they consider the task of deciding where to offload tasks given a single UE moving through multiple cells with MECs of varying capacities.

Ma et al. [22] propose to build an efficient service handoff system across edge servers based on Docker container migration. Wang et al. [38] utilize a lightweight heuristic algorithm to reduce offloading task execution delay by jointly considering task information, small base station and user mobility information. Yu et al. [42] propose a dynamic algorithm for partial offloading based on short-term mobility prediction to minimize energy consumption while satisfying delay requirements.

To summarize, we find the existing MEC-aware handoff works **fall short** in several respects:

- (1) Most solutions do not provide realistic algorithms that take the X2 application protocol into consideration [5, 13, 20, 22, 27, 38, 42].
- (2) Some of the proposed algorithms have computational complexities larger than $O(n^2)$ [38, 42]. In cases with larger numbers of UEs, the complexity becomes problematic for real time operation.
- (3) Related proposals tend to require additional message transmissions between UE and base station to collect information for the handoff algorithm. This overhead impacts system performance through additional link transmissions and information collection delay.
- (4) Related works do not inform their simulations with empirical 5G measurements thus making their interpretations less reliable [5, 13, 20, 22, 27, 38, 42].
- (5) Most related works lack detailed simulations. The related works conduct only numerical modeling or simplified simulations without millisecond granularity and packet-level/multilayer detail (e.g., from physical to application layer), thus excluding some detailed dynamics only visible with such simulations [5, 13, 20, 27, 38, 42].

We address these shortfalls by proposing an easily-deployable handoff algorithm with minimum overhead. We code the key metric values collected from a real 5G testbed together with Comp-HO algorithm into an open source custom MEC-enabled NS-3 simulator following base station protocol standard and perform a packet-level simulation.

²In this work, we define user experienced delay as the time between when user sends a request and receives a reply.

Table 1: Notation Table

и	UE u
m	MEC m
${\cal S}$	Serving base station and MEC server
${\mathcal T}$	Targeting base station and MEC server
Q^{m}	Processing queue length in MEC m
D_u^m	Transmission delay from u to m
S_u^m	Signal quality for u received from the cell with m
R_u	UE measurement: $\langle S_u^m, A \rangle$
θ	RSRQ threshold
δ	Handoff offset

3 PROBLEM FORMULATION

This section formulates the *MEC HO* problem in the 5G context with MEC servers co-located with 5G base stations. The formulation focuses on optimizing the performance of MEC-assisted UE applications with respect to experienced delay. Also we assume the transmission delay between a MEC server and its co-located base station is negligible (as also shown in the 5G testbed measurements, see Figure 1). We also note that there are many MEC location deployment schemes, e.g., co-located BS and MEC or several BSs sharing MEC co-located with an MME. Since comp-HO can be easily generalized to work these with different deployments (by considering MEC load only for HOs between BSs with different MECs), we only focus on the co-located BS with MEC case.

We let $\mathcal{M} \triangleq \{1, 2, ..., M\}$ denote the M MEC servers and $\mathcal{U} = \{1, 2, ..., U\}$ the U mobile UEs in the system. Each MEC server has a fixed capacity c (maximum queue length) and is connected directly to a 5G base station via fixed-line Ethernet. The time horizon is discretized into slots of equal periods indexed by $t \in \mathbb{N}$. We note several important MEC server assumptions:

- Homogeneity of MEC servers: The MEC servers have the same data processing rates.
- (2) Job-level migration: The tasks of a job may execute on any given MEC server.
- (3) Task atomicity: A task cannot be split across MEC servers.
- (4) Non-preemptive task scheduling: A task being processed cannot be interrupted by any other task.

Let D_u^m denote the transmission delay from UE u to MEC m excluding the processing delay. In other words, D_u^m consists of the uplink and downlink delay. Normally, either the uplink or downlink has larger data packet sizes and thus dominates the transmission delay. For example, uplink delay dominates the transmission delay, since the uplink packets to MEC servers contain much more data than the downlink packets for most MAR offloading applications. To simplify the problem, we consider only the dominant direction of data transmission, D_u^m , which solely depends on the signal quality received by u from m, i.e., S_u^m :

$$D_u^m = \phi(S_u^m). \tag{1}$$

where the function $\phi: \mathbb{R}^1 \to \mathbb{R}^1$ is monotonically decreasing.

Let Q^m denote the processing queue length of MEC m at a point in time. We assume the change in user experienced delay during the handoff from MEC S to T depends on the difference in signal

```
Algorithm 1: Comp-HO algorithm
    parameter: S_u^m \leftarrow \text{RSRQ } u received from m Q^m \leftarrow \text{Max queuing time in } m
     UE Measurement
     thread ReportUeMeasurment(R_u):
         while u offloading to m do
              \mathbf{S}_{u}^{\mathbf{m}} \leftarrow S_{u}^{m} for all probeable MECs, \mathbf{m} \subseteq \mathcal{M}
  2
              R_u \leftarrow \langle S_u^{\mathbf{m}}, A \rangle
                                                     //A - > App info
               sendUeMeasurement(R_u)
     Handoff
     thread updateUeMeasurement(R<sub>11</sub>):
  5 | \mathbf{R_u} \leftarrow R_u from all connected UEs, \mathbf{u} \subseteq \mathcal{U}
     thread updateLoad(Qm):
         Q^m, A \leftarrow \langle Q^m, A \rangle from all nearby MECs, m \subseteq M
     thread Hand off u:
         if S_{u}^{S} < \theta then
                                             // \theta-> RSRQ threshold
              8
                   SendHoRequest(T)
 10
```

qualities and queue lengths as follows:

$$\triangle D_{u}^{\mathcal{S},\mathcal{T}} = f(\triangle S_{u}^{\mathcal{S},\mathcal{T}}, \triangle Q^{\mathcal{S},\mathcal{T}}), \tag{2}$$

where $\triangle X_u^{\mathcal{S},\mathcal{T}} = X_u^{\mathcal{T}} - X_u^{\mathcal{S}}$. The function f denotes that both signal strength and computational load are considered.

Let $a_u^{\mathcal{S},\mathcal{T}}$ indicate whether MEC \mathcal{S} hands off u to \mathcal{T} as follows:

$$a_u^{\mathcal{S},\mathcal{T}} = \begin{cases} 1 & \text{if } \mathcal{S} \text{ hands off task of } u \text{ to } \mathcal{T}, \\ 0 & \text{otherwise, including no handoff requests.} \end{cases}$$
 (3)

We can then formulate an optimization problem as follows,

min
$$\sum_{u \in \mathcal{U}} \Delta D_u^{\mathcal{S}, \mathcal{T}} a_u^{\mathcal{S}, \mathcal{T}}$$
s.t. $\mathcal{S}, \mathcal{T} \in \mathcal{M}, \mathcal{S} \neq \mathcal{T}$ (4)

Local vs Global: The problem aims at optimizing user experience by minimizing overall user experienced delay for all UEs. The problem can be seen as a linear sum assignment problem which is also known as a minimum weight matching in bipartite graphs. Balanced assignment algorithms such as the Hungarian algorithm [9] can be used to solve this problem. Therefore, an optimal solution to the problem with global information is possible.

However, a global solution faces three challenges: 1) the potential for a single point of failure (the node making the global decisions), 2) the delay overhead caused by global information collection and decision dissemination can degrade system performance, and 3) most real-world standards require base stations to make their own hand-off decisions (thus modification of standards would be required). Therefore, we instead develop *Comp-HO* as a local optimization algorithm and leave global optimization for future work.

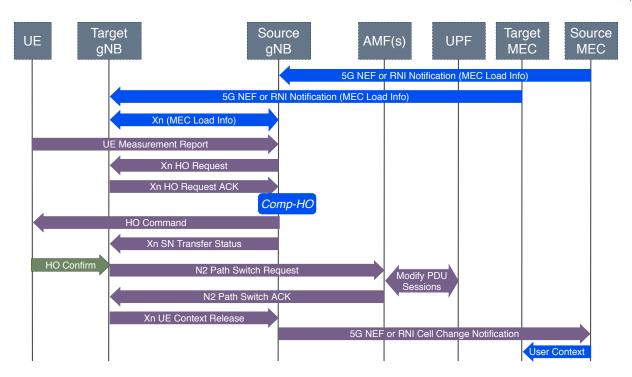


Figure 2: High-level Comp-HO flow diagram using 3GPP 5G and ETSI MEC [14] terminology. Comp-HO specific flows are marked in blue.

4 COMP-HO ALGORITHM

In this section, we first describe the Comp-HO handoff algorithm. Then we illustrate the integration of Comp-HO into the existing ESTI MEC standard information flows. Finally, we discuss the scalability of Comp-HO.

4.1 Algorithm

To perform a *MEC HO*, each UE sends measurement reports to the serving base station. The serving base station runs the handoff algorithm and decides when to initialize a handoff to a target base station. We first describe the UE measurement and then the handoff algorithm. Table 1 summarizes the related notations. Each UE collects the signal qualities (RSRQ) from all nearby base stations and sends them to the serving base station together with the information of its MEC-assisted applications in the report as shown on line 4 in Algorithm 1.

In parallel, each base station collects UE measurements from connected UEs and load information from nearby MEC servers, respectively (line 5 to 6). The load information includes metadata of MEC-assisted applications running in each MEC server and the processing queue lengths. $F(S_u^m, Q^m)$ denotes the weighted sum of signal quality and queue length to perform the optimization. To keep the complexity low, F() follows a linear form: $F(S_u^m, Q^m) = w_s * S_u^m - w_q * Q^m$. In the simulation, we iteratively tried different sets of weights and offsets to optimize performance. The base station starts the handoff of a UE if its RSRQ fails to meet the threshold (line 7). Utilizing the collected load information and signal quality values, the base station selects the target base station

and MEC server and sends the handoff request (line 8-10). The offset metric, δ , is introduced to avoid the ping-pong effect (line 9).

4.2 Flow Diagram

In terms of interaction between 5G network elements, Figure 2 illustrates the high-level message flows between such elements just before and during a *Comp-HO*. The figure uses the terminology and interfaces following standard 3GPP 5G and ESTI MEC [14], except that the Network Exposure Function (NEF) or Radio Network Interface (RNI) should allow passing load information back to the gNB/ng-eNB to be used in the *Comp-HO* algorithm. These NEF/RNI interfaces are currently designed to provide radio network information such as cell change notifications to the MEC system for user context (e.g., virtual machine) migration between MEC servers.

Also, we note that the analogous flows with a 5G radio access network and LTE EPC (like in testbed network) would be only slightly different in functions and terminologies. Overall, the flows illustrate how *Comp-HO* could potentially function after integrated into current mobile standards.

4.3 Scalability

The local algorithm is O(n) in time complexity where n is the number of sector UEs. This complexity is the same as the baselines (single signal indicator based handoff algorithms), so the algorithm scales in time at least as well as those. Whereas in terms of measurement messages, in our scenario the UE does not deal with MEC load info messages directly but these messages are transferred between the MECs and base stations over Xn links (as Figure 2 illustrates).

Thus the messaging complexity scales with the number of MECs and base stations rather than UEs, thus allowing good scaling as the number of UEs increase.

5 5G MEC MEASUREMENTS

5.1 5G MEC Testbed and MAR Prototype

Our 5G MEC testbed is a 5G micro-operator [24] network built by a joint national effort of academic and industrial partners. As shown in Figure 1, the testbed is composed of two parts: a 5G Core [1] and a base station. The core network is deployed in non-standalone (NSA) mode and the network functions (NFs) are implemented in Linux virtual machines located on servers in a single university server room. The MEC server is located in the same server room and connected via Ethernet to the core network switches. The base station is located on the roof of another building. The base station has two antennas for 5G and one for LTE, providing coverage over the campus area. The base station is connected to the core network via optical fiber. According to our measurements, the median RTT between the base station and MEC is 0.5 ms.

We develop a MAR prototype by running a custom Android client app on a Huawei Mate 30 Pro 5G smartphone (the UE) and a Linux server app on the MEC equipped with 8-core Xeon CPU, 16 GB memory and Quadro K2200 GPU. The client app captures camera frames at 10 frame rate (FPS). Then, it downscales the frames to 480×320 pixels and sends to the MEC.

The MEC receives the frames and uses YOLO [31] to perform object detection. The object detection result is composed of a set of bounding boxes of the detected objects as well as the object classes and detection confidences. Once the objects are detected, the result is sent back to the UE and rendered on the screen, annotating the objects from camera view.

Unfortunately, due to technical and licensing limitations, we cannot implement and deploy the *Comp-HO* algorithm into the base station. Instead, we conduct network measurements to observe the baseline performance of the MAR prototype system in Section 5.2. Then, in Section 6, we use these measurement results to inform the simulation parameters thus helping to mitigate the gap between a real-world 5G network and the NS-3 simulator.

5.2 Measurement Setup and Results

We record the transmission delays of the frames and detection results by timestamping during the sending and receiving on the UE and MEC. We connect the UE and the MEC via Ethernet beforehand to estimate the clock drift (within a confidence interval of <1 ms). We also select an off-peak time to conduct the measurements. The results, therefore, illustrate MAR performance with 5G MEC without non-MAR loads from other UEs.

We walk along a fixed route during the measurement and collect the results shown in Figure 3 and Table 2. Overall, the UE sends 540 frames to the MEC over a period of 54 seconds. The median frame (uplink) transmission delay and the result (downlink) transmission delay are 32.0 ms and 2.0 ms, respectively. This large difference is due to two primary reasons.

(1) Firstly, the uplink bandwidth is much smaller than the downlink. According to our measurements, the downlink throughput is about 360 Mbps while the uplink only 30 Mbps.

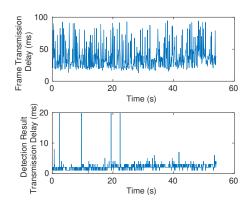


Figure 3: Frame transmission delay and detection result transmission delay between the UE and the MEC, measured with the MAR prototype in 5G MEC testbed.

Table 2: Measurement statistics of Figure 3. UE \rightarrow MEC: frame transmission (uplink). MEC \rightarrow UE: detection result transmission (downlink).

	$UE \rightarrow MEC$	$MEC \rightarrow UE$
Median delay (ms)	32.00	2.00
Jitter (ms)	9.5	0
Packet loss (%)	0.06	0.52

(2) Secondly, a frame has much larger size than its detection result. Specifically, a frames after compression is typically a few kilobytes while the detection result is simply plain text only typically only tens of bytes.

The combination of these two factors contributes to a higher uplink transmission delay than downlink transmission delay.

6 SIMULATION

6.1 NS-3 Simulator and Setup

We next perform simulations to estimate the performance of the *Comp-HO* algorithm at scale. The simulation setup aligns with some important metrics and results taken from the 5G measurements such as the frequency, UE speed and lower bound of user experienced delay.

Simulator: We modify the existing LTE module³ of NS-3 so that each eNB is co-located with three MEC servers (one for each sector) and any UE data packets are forwarded to a MEC server rather than to a packet gateway (Figure 4). We also integrate the *Comp-HO* algorithm into the LTE module. Each MEC server contains several queues that each process packets at a fixed rate. For a given MEC server, the processing queue length reported to the eNB (and thus to the handoff algorithm) is the minimum length out of these queues.

 $^{^3 \}mbox{We}$ use the LTE module rather than a 5G NS-3 module because handoff support in such modules [25, 29] is still limited.

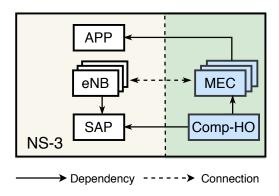


Figure 4: Custom NS-3 simulator

Transmission Setup: For simplicity, each UE sends UDP packets at a fixed FPS (20 Hz) but with random starting times to avoid synchronization. As mentioned, the packets are forwarded by the eNB to the corresponding MEC for processing and after processing sent back to the UE. If after processing, the UE is no longer being served by the eNB that received the packet originally (for example because of handoff) the packet is discarded. We denote this type of packet loss as MEC-mobility packet loss ⁴. For performance monitoring, the round-trip packet delay (which we refer to as the UE experienced delay) for each packet is recorded at the UE. We also track the number of handoffs.

Variations: We perform both basic simulations (with the noted parameters) and several simulation variations as follows.

- (1) *Handoff rate*: We vary the UE speeds to alter the effective handoff rate.
- (2) FPS: We vary the FPS to uniformly alter the MEC loads (given the fixed processing rate),
- (3) *Mobility:* We use two different UE mobility models (one with a center bias and one without) to illustrate the impact of different spatial UE distributions (e.g., a central crowd of UEs).

For further clarification on the mobility models, we use random waypoint as the baseline model as random waypoint has a center bias [6] thus allowing a higher central user density and heterogeneous load across MECs. Such heterogeneity naturally helps illustrate the benefits of the *Comp-HO* algorithm. However, for reference, we also use a Gauss-Markov mobility model which does not have such a center bias [8], and thus has a more homogeneous user distribution and load across MECs. We also perform each simulation three times (with different random seeds) to ensure performance differences are not due to random variation.

Benchmark: As pointed out in section 2, most MEC-aware hand-off proposals do not take the X2 protocol into consideration [5, 13, 20, 22, 27, 38, 42] and thus are incompatible with multilayer simulators such as NS-3, or have significant complexity (larger than

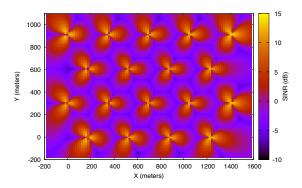


Figure 5: NS-3 simulation network layout map with SINR

Table 3: NS-3 simulation parameters

Parameter	Value		
Number of base stations	18 Tri-sector eNBs		
Layout	Hexagonal		
Intersite Distance	350 m		
Center Frequency	3.55 Ghz (band 22)		
Bandwidth	20 Mhz		
Path Loss Model	ITU-R P.1411 LoS [17]		
Height	eNB: 45 m, UE: 1.5 m		
Number UEs	50		
UE Mobility	Random Waypoint		
UE Velocity	2 m/s (7.2 km/h)		
Queues per MEC Server	1		
Simulation Area	1800x1300 m		
Simulation Time	30 s		

Table 4: Benchmark Parameters

Algorithm	Parameters
A2-A4-RSRQ	ServingRsrqThreshold=30 [2]
	NeighbourRsrqOffset=1 [2]
A3-RSRP	TimeToTrigger=256 (ms) [2]
	Hysteresis=3 (dB) [2]

 $O(n^2)$ [38, 42]) and thus may require further improvements for practical use. Additionally, to the best of our knowledge there are no available open source implementations of MEC-aware handoff algorithms (besides ours [4]). Therefore, we leave comparisons with other MEC-aware handoff algorithms for future work when more implementations are available.

Instead, we compare the *Comp-HO* algorithm to two existing LTE handoff algorithms (*A2-A4-RSRQ* and *A3-RSRP*) and to a scenario with no handoffs (NoHO). The *A2-A4-RSRQ* and *A3-RSRP* algorithms are based on the A2, A4 and A3 control events defined

⁴This loss is a challenge for MEC systems with LTE as the system cannot easily impact LTE traffic routing, thus forcing somewhat complex solutions. Luckily, 5G has flexible user plane functions that the external MEC system can change, thus reducing the issue. In any case, even removing this type of loss completely does not qualitatively change our delay results.

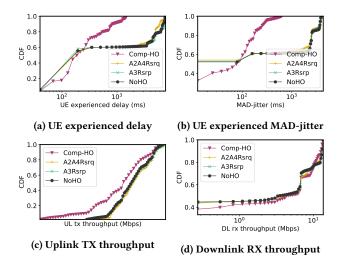


Figure 6: Performance improvements (x-axes on log scale)

by 3GPP standard [3]. The *A2-A4-RSRQ* algorithm triggers a hand-off when the serving cell's RSRQ drops below a threshold (A2) and a neighboring cell's RSRQ rises above an offset (A4). The *A3-RSRP* algorithm triggers a handoff when the serving cell's Reference Signal Received Power (RSRP) drops below the RSRP of a neighboring cell. Both algorithms are already part of the LTE module of NS-3 [28]. The other simulation parameters are summarised in Table 3. The developed NS-3 code is available at [4].

6.2 Transmission Performance

Improvement: Figure 6a and Figure 6b illustrate the cumulative distributions of UE experienced delay and jitter for all packets from UEs. We note that since the UE experienced delay distributions are very broad, we use a robust median-absolute deviation (MAD) version of jitter. As shown, *Comp-HO* improves the UE experienced delay with considerable improvements at the tail of the distribution. While, for jitter, *Comp-HO* essentially smooths out the distribution, thus removing the large bifurcation seen in the other algorithms. *Comp-HO* does this by better distributing the UEs across MECs and thus avoiding the issue of lucky UEs that happen to be in the areas with fewer competitors for channel and MEC resources.

Figure 8b and Figure 9b further support this conclusion by showing that with *Comp-HO* the MEC servers actually process more packets. Relatedly, the throughput performances in ?? show that with *Comp-HO* the UEs receive more processed packets (than the benchmarks) despite actually sending slightly fewer packets to be processed (than the benchmarks). The UEs with Comp-HO send slightly fewer packets packets because of a higher HO rate, which we discuss further in the next section.

To concentrate on the majority of the performance, in the rest of the paper we calculate the mean values excluding the effect of the outliers⁵. Overall, *Comp-HO* outperforms the benchamark algorithms with improvements in mean UE experienced delay of 73%-80% and downlink throughput of 3%-5%. Table 5 summarizes

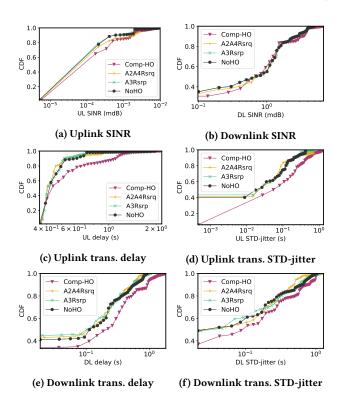


Figure 7: Trade-off (all x-axes on log scale). The transmission delays have excluded all processing delays.

Table 5: Performance improvements

	Сотр-НО	A2-A4-RSRQ	A3-RSRP	NoHO
Delay (ms)	214	788	1052	1096
DL TP (Mbps)	5.28	5.13	5.06	5.03
UL TP (Mbps)	6.68	to 6.74	6.73	6.73
MEC packets	27096	24889	24366	24433

the numerical results of UE experienced delay, downlink/uplink throughput, and MEC processed packets.

Trade-off: To illustrate the potential trade-off in using *Comp-HO*, Figure 7 illustrates the uplink and downlink SINR, transmission delay and jitter, therefore isolating transmission dynamics by removing processing dynamics. The results show that *Comp-HO* has somewhat lower SINR and higher transmission delay and jitter, as would be expected given that *Comp-HO* does not purely optimize SINR. However, the increases are relatively minor in comparison to the UE experienced delay and jitter decreases in the simulation. We also note that, as expected, the MEC-mobility packet loss of *Comp-HO* is 4.3% compared to 1.2% for *A2-A4-RSRQ*. This is the result of *Comp-HO* more eagerly switching base stations given high delay MECs. However, much of this loss could be avoided in ETSI MEC networks that include MEC assisted user context transfer as the transfer could include data (packets) waiting for processing.

 $^{^{5}}$ An outlier is more than one standard deviation from the overall mean

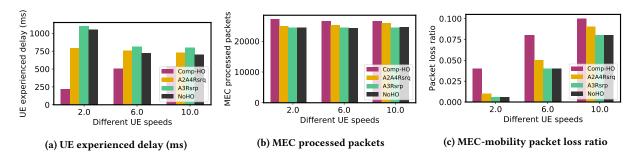


Figure 8: Performance with different UE speeds (2, 6, and 10 m/s)

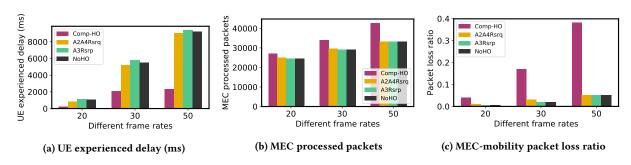


Figure 9: Performance with different frame rates (20, 30, and 50 fps)

Relatedly, the Comp-HO handoff rate is two to three times the benchmark rates depending on the specific scenario. This causes more service interruptions; however, when Comp-HO decides to handoff the service is in any case usually interrupted by MEC congestion. Thus the benefit outweighs the cost. Additionally, 5G techniques such as Dual Active Protocol Stack (DAPS) during handover should minimize the length of these handoff interruptions. The other significant consideration of higher handover rates is energy consumption, which we leave for future work.

6.3 Effects of Variations

Handoff rate: For different effective handoff rates (as varied through changes in UE speed), Figure 8c shows that *Comp-HO* does suffer higher MEC-mobility packet loss ratio (PLR) comparing with other algorithms with slower UEs. Meanwhile, Figure 8a and Figure 8b illustrate that higher handoff rates do decrease the performance gap between baseline algorithms and *Comp-HO*. However, *Comp-HO* still outperforms all baselines in all speeds for delay and MEC load assignment performance. The decreasing gaps are actually partly an artifact of the increase in MEC-mobility packet loss since more long-delayed packets are lost and not included in the delay distributions.

FPS: Figure 9 illustrate the performance comparison when using different frame rates (FPS), in other words different UE sending rates, within a common FPS range of MAR applications, i.e., 20 Hz, 30 Hz, and 50 Hz, respectively. The results show that *Comp-HO* improves user experienced delay significantly in all FPS, though the degree of improvement varies likely due to an interplay of factors. Due to the focus on both MEC load and signal strength, *Comp-HO*

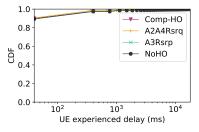


Figure 10: UE experienced delay with Gauss-Markov mobility model.

encounters higher packet loss ratio in higher frame rates scenario. Nevertheless, the result clearly shows that *Comp-HO* becomes more superior in terms of MEC load balancing and UE experienced delay while frame rate increases, indicating the benefit and necessity of *Comp-HO* for compute-intensive application offloading.

Mobility Models: For the two different mobility models, Figures 6a and 10 illustrate the UE experienced delay distributions of random waypoint and Gauss-Markov model respectively. As described before, with the baseline center-biased random waypoint model (with more heterogeneous user density and MEC loads) the *Comp-HO* algorithm provides significant gains. However, with the Gauss-Markov model (with more homogeneous user density and MEC loads) the *Comp-HO* algorithm provides very little benefit because there is less congestion at MEC servers and thus less need for load redistribution.

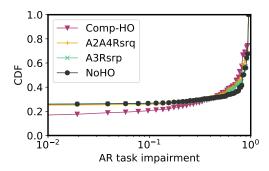


Figure 11: AR task impairment distributions. An impairment score of one represents the maximum empirical performance while a score of zero represents the minimum, thus a higher score is better.

6.4 Quality of Experience

To illustrate the potential impact on user QoE, we utilise an existing AR task impairment model [18] to transform the delay distribution values into AR task impairment scores. Specifically, an impairment score is a normalised score that quantifies the reduction in performance (from an empirical maximum) on a task or game (e.g., game score), specifically in this case, collaborate assembly of a virtual AR object like in Minecraft Earth (though also with physical object detection in our case). In other words, an impairment score of one represents the maximum empirical performance while a score of zero represents the minimum. Figure 11 illustrates the impairment score distributions for the different handoff algorithms. We find that *Comp-HO* significantly decreases the fraction of packets with full impairment (indicating the lowest performance level), thus suggesting that the delay improvements from *Comp-HO* should translate into actual QoE gains during these types of AR tasks.

6.5 Takeaway

The simulation results show that *Comp-HO* significantly improves the user experienced delay at the expense of a small increase in transmission delay (due to a decrease in signal strength). Additionally, *Comp-HO* is more robust to different UE speeds and outperforms benchmark algorithms in different FPS. Comparing with homogeneous user density and MEC loads distribution, *Comp-HO* provides more improvement of user experienced delay in heterogeneous counterpart.

7 DISCUSSION

In this work we proposed and evaluated *Comp-HO* as a simple, effective, and standard-friendly computational handoff algorithm. In the future though, both *Comp-HO* and our evaluation methods can still be improved in several aspects.

Implementation: In terms of our evaluation, firstly, the 5G tesbed uses commercial BSs and thus due to technical and legal reasons we could not implement *Comp-HO* directly into those BSs. 5G testbeds such as COSMOS [30] and POWDER [7] provide large scale open SDR networks, however they currently have limited access for other researchers and moreover it would be difficult to emulate

UE mobility without field test. Secondly, the custom MEC-enabled NS-3 simulator could include some features like the consideration of user context migration between MEC servers (e.g., the migration of docker containers or VMs containing the current app state [22, 44]) and a more realistic MAR model at the application level.

Algorithm: In terms of improvements to *Comp-HO* itself, the algorithm adds overheads such as collecting and sending the MEC load information. Thus, in cases with very few or very spatially homogeneous users, the algorithm will essentially act as a strongest SINR algorithm (since the MEC loads will be similar) but with an additional overhead and thus less efficient. Adding a threshold based on the number or distribution of UEs could remedy this by allowing base stations to fall back to traditional handoff algorithms during the mentioned conditions. Additionally, similar to having the NS-3 simulator consider the user context migration in the evaluation, *Comp-HO* could consider the user migration cost in the optimization directly (assuming differing migration costs between different MEC pairs).

Next step: In future work, we will integrate *Comp-HO* and related algorithms into open Software Defined Radio (SDR) base stations, therefore allowing evaluations in a real physical deployment. We also plan to further develop the MEC-enabled NS-3 simulator to allow more comprehensive simulations. We will develop algorithms that consider different types of interference in 5G networks and the tradeoff between the handoff rates and the level of interference in the network. In addition, we will develop a predictive capability algorithm allowing us to predict better mobility so resources can be pre-allocated before a handoff event even occurs.

8 CONCLUSION

The combination of 5G and MEC allows novel services and improved experiences in areas like MAR and virtual reality among others. Towards this goal and with a focus on mobility, this work studied the issue of *MEC HO* in 5G and proposed a handoff algorithm, *Comp-HO*, that considers both network signal strength and nearby MEC server load. We conducted the 5G MEC testbed measurements with a MAR prototype and utilized the collected results to set up large scale simulations with a custom MEC-enabled NS-3 simulator. The simulation results illustrated that *Comp-HO* improves the end-to-end delay compared to benchmark handoff algorithms in MEC scenarios. As far as we know, this is one of the first efforts to optimize *MEC HO* for 5G.

REFERENCES

 3GPP. 2019. Release description; Release 15. Technical report (TR) 21.915. 3rd Generation Partnership Project (3GPP). V15.0.0.

- [2] 3GPP. 2020. Evolved Universal Terrestrial Radio Access (E-UTRA); Radio Resource Control (RRC); Protocol specification. Technical specification 36.331. 3rd Generation Partnership Project (3GPP). V16.1.1.
- [3] Evolved Universal Terrestrial Radio Access. 2013. Radio resource control (RRC). Protocol specification (Release 10) (2013).
- [4] Authors. 2020. Source to be released on publication.
- [5] Fani Basic, Atakan Aral, and Ivona Brandic. 2019. Fuzzy handoff control in edge offloading. In ICFC'19.
- [6] Christian Bettstetter, Hannes Hartenstein, and Xavier Pérez-Costa. 2004. Stochastic properties of the random waypoint mobility model. Wireless Networks (2004).
- [7] Joe Breen, Andrew Buffmire, Jonathon Duerig, Kevin Dutt, Eric Eide, Mike Hibler, David Johnson, Sneha Kumar Kasera, Earl Lewis, Dustin Maas, et al. 2020. POWDER: Platform for Open Wireless Data-driven Experimental Research. In WINTECH'20.
- [8] Dan Broyles, Abdul Jabbar, James PG Sterbenz, et al. 2010. Design and analysis of a 3–D gauss-markov mobility model for highly-dynamic airborne networks. In ITC'10
- [9] Derek Bruff. 2005. The assignment problem and the hungarian method. Notes for Math (2005).
- [10] P Carter and MA Beach. 1995. Evaluation of handover mechanisms in shadowed low earth orbit land mobile satellite systems. *International Journal of Satellite* Communications (1995).
- [11] Dimitris Chatzopoulos and Pan Hui. 2016. Readme: A real-time recommendation system for mobile augmented reality ecosystems. In ACM MM'16.
- [12] Mostafa Zaman Chowdhury, Won Ryu, Eunjun Rhee, and Yeong Min Jang. 2009. Handover between macrocell and femtocell for UMTS based networks. In ICACT'09.
- [13] Mustafa Emara, Miltiades C Filippou, and Dario Sabella. 2018. MEC-aware cell association for 5G heterogeneous networks. In WCNC'18.
- [14] European Telecommunications Standards Institute (ETSI). 2019. Multi-access Edge Computing (MEC); Application Mobility Service API. Recommendation 012 V2.1.1.
- [15] Ilija Hadžić, Yoshihisa Abe, and Hans C Woithe. 2017. Edge computing in the ePC: A reality check. In SEC'17.
- [16] J. Han and B. Wu. 2010. Handover in the 3GPP long term evolution (LTE) systems. In 2010 Global Mobile Congress.
- [17] ITU-T. 2019. Propagation data and prediction methods for the planning of short-range outdoor radio communication systems and radio local area networks in the frequency range 300 MHz to 100 GHz. Recommendation P.1411. International Telecommunication Union.
- [18] B. Krogfoss, J. Duran, P. Perez, and J. Bouwen. 2020. Quantifying the Value of 5G and Edge Cloud on QoE for AR/VR. In QoMEX'20.
- [19] Janne Kurjenniemi, Tero Henttonen, and Jorma Kaikkonen. 2008. Suitability of RSRQ measurement for quality based inter-frequency handover in LTE. In ISWCS'08.
- [20] Tong Li, Chathura Sarathchandra Magurawalage, Kezhi Wang, Ke Xu, Kun Yang, and Haiyang Wang. 2017. On efficient offloading control in cloud radio access network with mobile edge computing. In ICDCS'17.
- [21] Luyang Liu, Hongyu Li, and Marco Gruteser. 2019. Edge assisted real-time object detection for mobile augmented reality. In MobiCom'19.
- [22] Lele Ma, Shanhe Yi, and Qun Li. 2017. Efficient service handoff across edge servers via docker container migration. In SEC'17.
- [23] Yuyi Mao, Jun Zhang, and Khaled B Letaief. 2017. Joint task offloading scheduling and transmit power allocation for mobile-edge computing systems. In WCNC'17.
- [24] Marja Matinmikko, Matti Latva-Aho, Petri Ahokangas, Seppo Yrjölä, and Timo Koivumäki. 2017. Micro operators to boost local service delivery in 5G. Wireless Personal Communications (2017).
- [25] M. Mezzavilla, M. Zhang, M. Polese, R. Ford, S. Dutta, S. Rangan, and M. Zorzi. 2018. End-to-End Simulation of 5G mmWave Networks. *IEEE Communications Surveys & Tutorials* (2018).
- [26] J-M Moon and D-H Cho. 2009. Efficient handoff algorithm for inbound mobility in hierarchical macro/femto cell networks. IEEE Communications Letters (2009).
- [27] Wahida Nasrin and Jiang Xie. 2019. A Joint Handoff and Offloading Decision Algorithm for Mobile Edge Computing (MEC). In GLOBECOM'19.
- [28] NS-3. 2020. Design Documentation Model Library. https://www.nsnam.org/docs/models/html/lte-design.html#handover.
- [29] Natale Patriciello, Sandra Lagen, Biljana Bojovic, and Lorenza Giupponi. 2019. An E2E simulator for 5G NR networks. Simulation Modelling Practice and Theory (2019).
- [30] Dipankar Raychaudhuri, Ivan Seskar, Gil Zussman, Thanasis Korakis, Dan Kilper, Tingjun Chen, Jakub Kolodziejski, Michael Sherman, Zoran Kostic, Xiaoxiong Gu, et al. 2020. Challenge: COSMOS: A city-scale programmable testbed for experimentation with advanced wireless. In MobiCom'20.

[31] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. 2016. You Only Look Once: Unified, Real-Time Object Detection. In CVPR'16.

- [32] Debashis Saha, Amitava Mukherjee, Iti Saha Misra, and Mohuya Chakraborty. 2004. Mobility support in IP: a survey of related protocols. IEEE network (2004).
- [33] Stefania Sardellitti, Gesualdo Scutari, and Sergio Barbarossa. 2015. Joint optimization of radio and computational resources for multicell mobile-edge computing. IEEE Transactions on Signal and Information Processing over Networks (2015).
- [34] M. Satyanarayanan, P. Bahl, R. Caceres, and N. Davies. 2009. The Case for VM-Based Cloudlets in Mobile Computing. IEEE Pervasive Computing (2009).
- [35] Wei Shen and Qing-An Zeng. 2006. A novel decision strategy of vertical handoff in overlay wireless networks. In NCA'06.
- [36] Kari Sipila, Mika Jasberg, Jaana Laiho-Steffens, and Achim Wacker. 1999. Soft handover gains in a fast power controlled WCDMA uplink. In VTC'99.
- [37] Nishith D Tripathi, Jeffrey H Reed, and Hugh F VanLandinoham. 1998. Handoff in cellular systems. IEEE personal communications (1998).
- [38] Zi Wang, Zhiwei Zhao, Geyong Min, Xinyuan Huang, Qiang Ni, and Rong Wang. 2018. User mobility aware task assignment for mobile edge computing. Future Generation Computer Systems (2018).
- [39] Peng Xu, Xuming Fang, Rong He, and Zheng Xiang. 2013. An efficient handoff algorithm based on received signal strength and wireless transmission loss in hierarchical cell networks. *Telecommunication Systems* (2013).
- [40] Peng Xu, Xuming Fang, Jun Yang, and Yaping Cui. 2010. A user's state and SINR-based handoff algorithm in hierarchical cell networks. In WiCOM'10.
- [41] Candy Yiu, Yujian Zhang, Mo-Han Fong, and Yuefeng Peng. 2015. User equipment and methods for handover enhancement using reference signal received quality (RSRQ). US Patent 9,130,688.
- [42] Fangxiaoqi Yu, Haopeng Chen, and Jinqing Xu. 2018. DMPO: Dynamic mobilityaware partial offloading in mobile edge computing. Future Generation Computer Systems (2018).
- [43] Zikang Yuan, Dongfu Zhu, Cheng Chi, Jinhui Tang, Chunyuan Liao, and Xin Yang. 2019. Visual-Inertial State Estimation with Pre-integration Correction for Robust Mobile Augmented Reality. In ACM MM'19.
- [44] Aleksandr Zavodovski, Nitinder Mohan, Suzan Bayhan, Walter Wong, and Jussi Kangasharju. 2018. Icon: Intelligent container overlays. In HotNets' 18.
- [45] Wenxiao Zhang, Bo Han, and Pan Hui. 2018. Jaguar: Low latency mobile augmented reality with flexible tracking. In ACM MM'18.
- [46] Ziyue Zhang, Jie Gong, Xiang Chen, and Terng-Yin Hsu. 2020. Reinforcement Learning Based Computation-aware Mobility Management in Ultra Dense Networks. Journal of Internet Technology 21, 6 (2020), 1785–1794.