

Gaming on the Edge: Performance Issues of Distributed Online Gaming

D. Olliaro¹, V. Mancuso², P. Castagno³, M. Sereno^{3,4}, M. Ajmone Marsan²

¹ Ca' Foscari University of Venice, Venice, Italy

² IMDEA Networks Institute, Leganes, Spain

³ Università di Torino, Torino, Italy

⁴ Consorzio Nazionale Interuniversitario per le Telecomunicazioni (CNIT), Italy

Abstract—We study the performance of online games played over a platform that implements gaming as a service (GaaS) in a mobile network slice that hosts concatenated virtual network functions (VNFs) at the edge. The distributed gaming architecture is based on edge computing facilities, whose utilization must be carefully planned and managed, so as to satisfy the stringent performance requirements of game applications. The game manager must consider the latency between players and edge server VNFs, the capacity and load of edge servers, and the latency between edge servers used by interacting players. This calls for a careful choice about the allocation of players to edge server VNFs, aiming at extremely low latency in interactions resulting from player's commands. We develop an analytical model, which we validate with experiments in the wild, and show that, under several combinations of system parameters, deploying gaming VNFs at the edge can deliver better performance with respect to cloud gaming, in spite of the complexities arising from the distribution of gaming VNFs over edge servers. Our analytical model provides a useful tool for edge gaming systems performance prediction, thus supporting the management of GaaS applications.

Index Terms—Online gaming; GaaS; Performance model.

I. INTRODUCTION

In the early days, video games demanded expensive specialized hardware consoles with advanced graphics processing capabilities. Network connections were primarily for player interactions, limited to those not on the same gaming console. Today, the landscape has shifted dramatically. People now seek to indulge in their favorite games using their smartphones anytime, anywhere. Games have evolved into services delivered via ad hoc virtual network functions (VNFs), following the Gaming as a Service (GaaS) paradigm. Game logic and video processing occur within a VNF housed in a computing infrastructure, typically a cloud data center [1]. Players transmit their game commands to the cloud, where the VNF amalgamates them with others' commands to construct the video scene, which is then encoded and streamed back to players. RTP streams via UDP ports facilitate bidirectional communication. Additionally, the game VNF in the cloud must assess the

player's connection quality to determine the appropriate video quality. Consequently, the player's device is tasked solely with collecting commands, transmitting them along with connection feedback to the VNF, and decoding and displaying the received video stream. Only network customers with broadband access and low latency to game VNFs can properly experience online gaming. Hence, considering the increasing relevance of mobile gaming, and the evolutionary trend of radio access networks (RANs), we can expect that game VNFs will soon move from cloud to edge networks. Hence, we propose and study a GaaS framework built on VNFs residing on edge networks.

Processing player's commands with a VNF residing on one of the edge computing facilities, and generating the corresponding video stream at the edge, has obvious advantages in terms of response delay and network resource utilization. However, one of the most relevant issues is the allocation of players to VNFs, as players that connect to VNFs residing in different computing facilities could experience consistency issues. Since we cannot expect that all interacting players will connect to the same VNF, it will be necessary to transfer relevant commands from one VNF to the other.

We look at those issues considering a slice of a cellular network. When players join from distant locations, the slice management, and operation (MANO) VNF must map players onto edge VNFs, considering latency between edge VNFs and players, load, and processing capacity of edge VNFs, as well as the latency between VNFs that serve interacting players. In the following, we use "edge server" and "game VNF" interchangeably, as well as "gamer" and "player".

Our primary contribution is the development of *the first stochastic model* for the evaluation of edge gaming *success probability*, i.e., the probability that player's commands be incorporated in the rendered game within a deadline compatible with the players' reaction time. Additional contributions are: *i)* the comparison of predictions obtained with our stochastic model with results of experiments that validate the simplifying model assumptions; *ii)* the assessment of the impact of the main system parameters on edge gaming success probability, using data derived from measurements of real gaming platforms; *iii)* a comparison of success probabilities achievable with edge or cloud gaming approaches.

This work has been supported by the Project AEON-CPS (TS1-063000-2021-38), funded by the Ministry of Economic Affairs and Digital Transformation and the European Union NextGeneration-EU in the framework of the Spanish Recovery, Transformation and Resilience Plan and by the RESTART Program, financed by the Italian government with the resources of the Italian Recovery, Transformation and Resilience Plan – Mission 4, Component 2, Investment 1.3, theme 14 "Telecommunications of the future" (PE00000001 - program "RESTART", CUP D93C22000910001, projects R4R and ITA NTN).

II. RELATED WORK

The concept of cloud gaming emerged around 2009 and brought to the online gaming domain the streaming technologies that Spotify and Netflix introduced in music and video domains about a decade earlier. However, cloud gaming designers soon realized that the quality of experience (QoE) requirements for online gaming are quite different from those for music and video. This led the ITU (International Telecommunications Union) Study Group 12, focusing on *Performance, QoS and QoE* [2] to work on the identification of factors that influence online gaming QoE and on approaches to predict, plan, and manage QoE [3]. Building on this initial activity, ITU SG 12 produced three recommendations: G.1032, P.809, and G.1072 [4]–[6]. Recommendation G.1032 identifies factors that influence cloud gaming QoE, while Recommendation P.809 identifies approaches for subjective evaluation of gaming QoE. Finally, Recommendation G.1072 describes QoE models for cloud gaming, accounting for the effect of underlying networks. Our work goes in the direction of Recommendation G.1072, except we develop a probabilistic approach for the estimation of the gaming QoS (which is not present in G.1072).

In [7] the authors define a methodology for the assessment of the user-perceived quality of cloud gaming systems and conduct extensive experiments. Other recent papers report measurement studies of cloud gaming applications [8]–[10], and the very recent paper [11] looks at the specific case of gamers connected through a cellular network. These works identify the main network characteristics of cloud gaming and look at achievable QoS metrics, from which they infer the most relevant drivers of user’s QoE. Measurements reported in these studies determined the choice of parameters we use later in this paper for the performance analysis of edge gaming.

It is interesting to note that all the above-mentioned works consider *cloud* gaming systems, i.e., systems where game control and graphics engines that generate video streams for gamers reside in the cloud. Instead, Kim *et al.* [12] propose to separate game control, which remains in the cloud, from rendering, that is brought to the edge, hence closer to the end user, with the advantage of reduced latency in video delivery, and lower impact on network resources, especially in the case of ultra high definition video. We bring this concept one step forward, allocating both game control and rendering engines on the edge, with additional benefits in terms of latency and resource consumption, at the cost resulting from coordinating edge computing facilities used by interacting gamers.

In [13], gaming is mentioned as a relevant use case for the IETF Reliable and Available Wireless (RAW) architecture. Moreover, in [14] the authors propose an emulation framework to investigate QoE in a real-time video streaming scenario.

Related studies explore architectural propositions for online games, relevant to our work. For instance, [15] proposes a distributed engine for online gaming. The paper [16] discusses edge computing’s role in addressing latency and network resource efficiency in online gaming. Unlike our work, these

proposals lack analytical insights for game service providers.

Some studies propose architectures for online gaming that prioritize the economical utilization of resources, encompassing bandwidth, computational resources, and energy. E.g., the authors of [17] leverage software migrations to optimize energy consumption, and [18]–[21] propose efficient resource allocation for massive multiplayer online gaming with latency constraints. Our work is complementary to resource allocation, because we analyze the compound importance of server characteristics, given the amount of allocated resources.

Demanding QoS prerequisites of massive multiplayer online gaming with strict latency constraints call for optimal and adaptable resource management of both computational and network resources. This particular challenge can be addressed through VNF optimization techniques similar to those proposed, for instance, in [22], [23]. For the specific case of online gaming, we offer a tool that permits to identify VNF requirements and VNF chaining to achieve high QoS.

III. SYSTEM MODEL AND GAME REQUIREMENTS

We consider a mobile network that supports a multiparty distributed gaming application with strict latency requirements. Gaming engines are located on the edge, spread across multiple game VNFs, with the goal of minimizing the delay experienced by end users. We cast the required gaming architecture onto the standard 3GPP (3rd Generation Partnership Project) architecture for communication and computing.

A. Connect-compute architecture

In a 3GPP network there are four main network functions that we need to consider [24]:

- *User Equipment* (UE) is the equipment on top of which the gamer plays. We consider wireless-connected gamers.
- *(Radio) Access Network*, or (R)AN is the (wireless) infrastructure that connects user’s devices to the core network. It includes several facilities enabling communication. We consider wireless access networks.
- *User Plane Function* (UPF) is a vital element in 3GPP (5G) networks, as it is responsible for managing data traffic and handles tasks such as packet forwarding, traffic routing, and QoS enforcement.
- *Edge* is a computing facility under the control of the cellular network, that brings computation resources closer to gamers, improving response times. Edges are deployed as VNFs accessible through the user plane (via UPFs). Their proximity enhances real-time interaction.

As shown in Fig. 1, the interaction between the mentioned components determines a chaining of VNFs. RAN establishes wireless links for reliable connectivity, while UPFs manage data transmission and routing. Different edge computing facilities can be reached from different UPFs with different delays.

B. Application requirements

Latency requirements depend on the level of interaction between gamers. If two gamers are closely interacting (e.g., fighting against the same monster), it is necessary that the

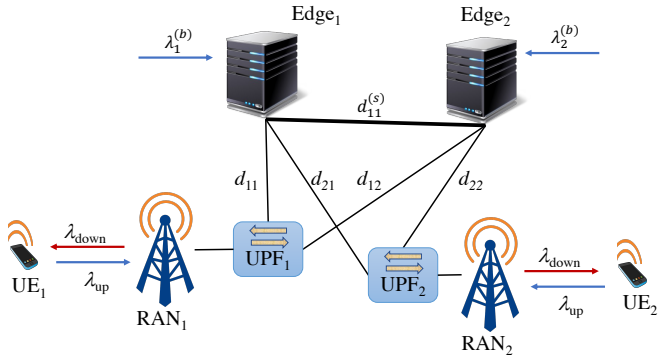


Fig. 1. Reference scenario

latency with which they observe each other's actions (and those of the monster) is as small and as equal as possible.

The most important QoS KPIs (key performance indicators) for online gaming are upstream and downstream data rates available between game server and user, upstream and downstream packet loss probabilities, and round trip latency. According to [8]–[10], typical requirements for immersive games are of the order of 1 Mb/s upstream (a small part for gamer's commands, the rest for feedback on network connection KPIs) and between 5 and 50 Mb/s downstream (for game video stream; of course, the downstream data rate depends on the gamer's screen size, and heavily impacts the quality of the video received by the gamer, hence the immersivity of the game). Packet loss probabilities are tolerated up to 10-30% in upstream (the KPI feedback is very redundant) and to a few percent in downstream. Latency requirements are of the order of 100 ms, which refer to the whole *game response delay* that includes all latency components due to network, processing, and buffering [7], [12].

IV. ANALYSIS AND PROBLEM FORMULATION

A. Reference scenario

The reference scenario for this work is depicted in Fig. 1, which is modeled as shown in Fig. 2. We consider a network with two or more base stations (BSs) where the (R)AN network function is accessed by UEs as a necessary step to access UPFs and two or more game VNFs. To simplify the description of the considered scenario, and given the geographically distributed nature of the service jointly with the extreme RAN densification process envisioned by 5G and beyond architectures, we assume that within the gaming slice, each BS hosts a single gamer and each UPF serves a single BS. Therefore, we use a single index to refer to a user/BS/UPF chain. This assumption can be easily removed—and was actually removed in the implementation of the model that will be described later in this section.

We assume that communicating between BS i and Edge j incurs three delay components: (i) a fixed delay d_{ij} ; (ii) a random component that accounts for user mobility, signal fluctuations, noise, scheduling, as well as other factors, including retransmissions due to losses; (iii) a variable queueing and processing time, which is modeled by means of an M/M/1 queue in uplink (from BS to Edge) and an M/M/1 queue in

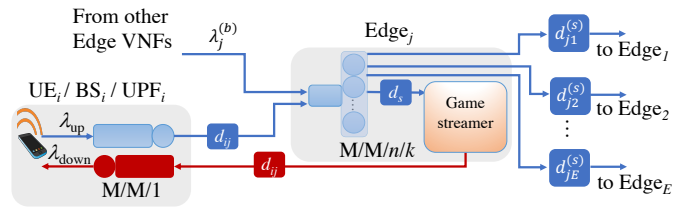


Fig. 2. Network model (detail for a generic gamer)

downlink (from Edge to BS). The sojourn time in the uplink (respectively downlink) queue associated with BS i and Edge j is a random variable denoted as U_{ij} (respectively D_{ij}).

An online gamer produces a flow of λ_{up} packets per second in uplink, and receives a multimedia stream of λ_{down} packets per second in downlink. For simplicity, we assume that the intensity of streams does not change over time. The stationary assumption is used because it allows to design and plan the service by considering periods of peak utilization. All uplink packets of a user are delivered by the BS to a game VNF, where they are processed by the edge server to which the user is associated. We model this uplink flow processing with a queue. We will primarily use a M/M/n/k model, $k \geq n$ (and we will present our modeling approach in this case), but we will also compare results against other models (and we will validate our assumptions with experiments). Hence, the game VNF is represented as an ensemble of n processors with a waiting room of $k-n$ packets. These parameters represent how resources are allocated within the network slice dedicated to online gaming. The time T_j spent by an uplink packet in Edge j is a random variable.

When a game VNF receives any instruction from a gamer that requires implementing changes in the downlink stream (e.g., a change of view, a movement that has to be reflected in the multimedia stream, etc.), the change is implemented with a fixed delay d_s , which represents the time needed to start producing the modified stream within the VNF.

B. Distributed gaming server model

Since gamers that are closely interacting are not necessarily close in the network, they may interface with different game VNFs and even different physical edge facilities. Thus, when two or more gamers interact, all streams they receive have to be adapted synchronously, with a tolerance of few tens of milliseconds. Hence, commands/moves issued by a gamer must be notified not only to the game VNF that manages the stream of that user but also to all other game VNFs involved.

We assume that gamers are not aware of the location of all game VNFs to be notified, while the serving game VNF of a gamer is. Therefore, when a command issued by a gamer is served by her Edge j , it is also forwarded to any Edge j' of any other gamer involved in the game scene, with a latency $d_{jj'}^{(s)}$ between serving game VNFs, which we model as a constant. This simplification is possible since edge servers are connected via high-speed networks (typically optical rings).

Guaranteeing a good level of interaction between gamers implies that the latency with which actions of one gamer are

visible to other gamers remain below a fixed threshold T_O (for example 100 ms) with sufficiently high probability (e.g., 0.8) and acceptable packet loss rate (a few percent) [7], [9].

C. Latency

In the simple case of two users, in the worst case, the close interaction of gamers involves two game VNFs and two BSs (each of which has a separate uplink and downlink processor) in addition to the VNFs located on gamers' terminals. In general, the worst case number of involved entities scales linearly with the number of interacting gamers. The resulting latency depends on communication and computing delays, and there are two groups of delays to evaluate:

- First, the *self-latency*, i.e., the delay with which a gamer i connected to Edge j sees her own actions reflected in her game's multimedia streaming, i.e.,

$$L_{ii}^{(j)} = d_{ij} + U_{ij} + T_j + d_s + D_{ij} + d_{ij}, \quad (1)$$

where the fixed latency between gamer i and Edge j is reasonably taken as symmetric in uplink and downlink.

- Second, the *cross-latency*, i.e., the latency between the moment a gamer i issues a command and any other gamer i' interacting with i and connected to Edge j' sees the effect of that command on her multimedia stream, i.e.,

$$L_{i'i'}^{(j')} = d_{ij} + U_{ij} + T_j + d_{j'j}^{(s)} + T_{j'} + d_s + D_{i'j'} + d_{i'j'}. \quad (2)$$

Let us consider for example the simple case of two gamers (indicated as g_1 and g_2) connected to two game VNFs (indicated as e_1 and e_2 , respectively). We need to account for four latency values: the round trip time between gamer g_1 and Edge e_1 , the latency of the path $g_1 \rightarrow e_1 \rightarrow e_2 \rightarrow g_2$, the latency of the symmetric path, $g_2 \rightarrow e_2 \rightarrow e_1 \rightarrow g_1$, and the round trip time between gamer g_2 and game VNF e_2 .

Assuming that it is possible to model with simple queues the delay introduced by individual system components (communication between gamers and game VNFs, and processing at game VNFs), and assuming independence among latency components, it is possible to compute probability distributions for relevant latency variables. E.g., it is possible to compute convolutions of distributions of additive delays or, equivalently, obtain the LST (Laplace–Stieltjes transform) of the delay pdf (probability distribution function) as the product of LSTs of pdfs of individual delay components.

D. Success probability

In addition to latency, if some queues have a finite capacity, we have to consider network losses. Hence, the success of multiparty communication depends on a combination of latency and loss. The success rate is the fraction of game control messages that are not lost and are processed on time to cause an observable change in the game multimedia streaming delivered to gamers in the same game scene, within a maximum latency T_O from message generation.

In our model, queues that can incur loss are the ones representing game VNFs, and the associated loss probability is

the probability that the buffer is full, which can be computed with standard queueing theory results.

Thus, once we know the loss probabilities at queues and the pdf of the latency, it is possible to compute the probability that the game interaction be completed within T_O for all gamers in the same game scene.

For the specific case of M/M/1, M/M/n/k and M/D/1-PS (M/D/1 with processor sharing) queues used in this paper, the following results hold:

- The LST of the sojourn time in an M/M/1 queue with arrival rate λ and service rate μ is

$$\hat{f}_{M/M/1}(s) = \frac{\mu - \lambda}{s + \mu - \lambda}. \quad (3)$$

- The LST of the sojourn time in an M/M/n/k queue with arrival rate λ and individual server service rate μ is

$$\hat{f}_{M/M/n/k}(s) = \frac{1}{1-p(k)} \frac{\mu}{s+\mu} \sum_{j=0}^{k-1} p(j) \left(\frac{n\mu}{s+n\mu} \right)^{[j-n+1]^+}, \quad (4)$$

with

$$p(j) = p(0) \frac{(\lambda/\mu)^j}{\min(j, n)! n^{[j-n]^+}}, \quad j = 0, \dots, k, \quad (5)$$

where $p(0)$ is found by solving $\sum_j p(j) = 1$ and $[x]^+$ indicates the max between 0 and x .

- The LST of the sojourn time in an M/D/1-PS queue with arrival rate λ and service time $1/\mu$ can be computed from (5.16) in [25], i.e.:

$$\hat{f}_{M/D/1-PS}(s) = \frac{\left(1 - \frac{\lambda}{\mu}\right) (\lambda + s)^2 e^{-\frac{\lambda+s}{\mu}}}{s^2 + \lambda \left(s + (\lambda + s) \left(1 - \frac{\lambda}{\mu}\right)\right) e^{-\frac{\lambda+s}{\mu}}}. \quad (6)$$

The random component of the latency due to the RAN activity is modeled through a random variable with limited variability. As suggested by experimental results reported in [26], [27], the RAN latency is close to the average most of the time, although values can vary in a relatively narrow interval. We therefore use a triangular distribution between latency values t_m and t_M , whose LST is as follows:

$$\hat{f}_{\text{RAN}}(s) = \left(\frac{2}{t_M - t_m} \cdot \frac{e^{-\frac{3t_M+t_m}{4}s} - e^{-\frac{t_M+3t_m}{4}s}}{s} \right)^2. \quad (7)$$

The above expression, with appropriate parameters t_m and t_M will be incorporated into network delays U_{ij} and D_{ij} .

The additive delay due to different path components has an LST equal to the product of the LSTs of the components. Of course, each LST related to processing must be tailored to the game VNF capacity in terms of number of servers, buffer capacity, server speed, and request arrival rate. Similarly, the LST relating to communication hops must be tailored to the link data rate, message size, and load. We omit here the specific expressions of latency LSTs because their derivation is straightforward.

The loss probability relative to arrivals at the game VNF is the probability that an arrival finds all k positions busy (k includes positions in service and in the waiting line). This probability, denoted by $p(k)$, equals the steady state probability of k positions busy, due to the assumption that the arrival process is Poisson, and is expressed as:

$$\pi_{\text{loss}} = p(k) = p(0) \frac{(\lambda/\mu)^k}{n! n^{j-n}}. \quad (8)$$

We use $F_{L_{ii'}}^{(j')}(t)$ to indicate the CDF (cumulative distribution function) of the cross-latency observed by game commands issued by a gamer i served by Edge j , and whose game control messages have to reach another gamer i' connected to a different Edge j' . $F_{L_{ii}}^{(j)}(t)$ indicates the self-latency. The probability of not exceeding the timeout is therefore simply computed as $F_{L_{ii'}}^{(j')}(T_O)$ and $F_{L_{ii}}^{(j)}(T_O)$, for cross- and self-latency, respectively. Notice that we assume that packet loss is negligible, including it in the RAN behavior. This assumption holds because gamer's transmissions use a modulation and coding scheme at which the bit error rate is very low, hence errors on a few bits can be recovered and the occurrence of error bursts can be obviated by means the HARQ protocol adopted at the BS. The impact of HARQ is accounted for in the model within the variable delay described by $\hat{f}_{\text{RAN}}(s)$.

Performance results depend on the intensity of flows that are offered to the modeled queues. In particular, apart from the above-mentioned λ_{up} and λ_{down} , we denote by $\lambda_j^{(b)}$ the background traffic of Edge j , generated by other gamers on the same server, playing in other game scenes. Notice that here we have neglected background traffic of communication links between gamers and game VNFs because we have assumed to use a network slice dedicated to a specific online game.

More complex and intertwined is the computation of the traffic offered to game VNFs. If we consider Edge j where multiparty gamers $G_j = \{g_1^{(j)}, g_2^{(j)}, \dots, g_{\ell_j}^{(j)}\}$ are connected and play also with remote groups of gamers $G_{j'} = \{g_1^{(j')}, g_2^{(j')}, \dots, g_{\ell_{j'}}^{(j')}\}$, the offered traffic is as follows:

$$\lambda_j = \left(\ell_j + \alpha \sum_{j' \neq j} \ell_{j'} (1 - \pi_{\text{loss}}^{(j')}) \right) \lambda_{\text{up}} + \lambda_j^{(b)}, \quad (9)$$

where $\pi_{\text{loss}}^{(j')}$ is the loss probability of Edge j' and α is the fraction of uplink messages that convey game control instructions instead of simple streaming feedback to the game server, hence they have to be forwarded to Edge j even when the gamers are served by Edge j' .

Eventually, the success probability of a gamer i in a multiparty game involving gamers connected to a set of E game VNFs can be expressed as the probability that all self and cross latency values originated at that gamer be below the timeout, after having discounted lost messages:

$$S_i = F_{L_{ii}}^{(j)}(T_O) \left(1 - \pi_{\text{loss}}^{(i)} \right) \cdot \prod_{j' \in E} \prod_{i' \in G_{j'} \setminus \{i\}} \left(F_{L_{ii'}}^{(j')}(T_O) \left(1 - \pi_{\text{loss}}^{(i)} \right) \left(1 - \pi_{\text{loss}}^{(i')} \right) \right). \quad (10)$$

The success of a multiparty game can be seen as the success of all interacting gamers:

$$S = \prod_{i \in \cup_j G_j} S_i. \quad (11)$$

which represents the fraction of players' actions that are incorporated into the video stream and visualized on gamer screens within the specified timeout T_O (i.e., those that are neither lost nor late).

V. PERFORMANCE EVALUATION

In this section we present the results of our analytical model, and we validate the model through experiments over the Internet. With experiments, we consider non-exponential distributions for the times between game commands generated by players, and for the VNF service times. In addition, we obviously remove independence assumptions.

Unless otherwise specified, we consider two edge servers (named e_1 and e_2) and a multitude of gamers. Two gamers (named g_1 and g_2) are modeled in detail, and individual performance metrics are collected for them. Other gamers are collectively modeled through the load they impose on game VNFs, which represents background traffic for the selected pair of gamers whose performance we study. Game VNFs are modeled as M/M/n/k queues with k taking values n (hence no buffering), $2n$ (hence the buffer space is as large as the number of servers n), or ∞ (hence an M/M/n queue). We assume $n=5$ at e_1 and $n=10$ at e_2 for the considered GaaS slice. We also consider an M/D/1-PS queue (where PS stands for Processor Sharing) at the edge servers, to approximate the analysis of realistic systems that implement round-robin execution of tasks (i.e., portions of jobs) rather than sequential processing of whole jobs.

The two gamers generate 1 Mb/s in uplink toward the game VNF to which they are connected, that is 1 000 packets per second, of 1 000 bits each; 10 of those packets correspond to game commands (1%), the rest is QoS feedback. Game VNFs generate 30 Mb/s in downlink toward each gamer (3 000 packets per second, each of 10 000 bits) to provide gamers with their game video stream. The time to prepare the video stream is assumed to be $d_s=10$ ms.

The values we use to parametrize our model were selected so as to reflect what was measured over popular online game platforms by independent studies, e.g., [8], [10]. We further assume that a single thread of a single processor on a server is able, on average, to serve a gamer. Thus, we fix $\mu=1\,000$ packets/s for the purposes of this numerical evaluation. The uplink and downlink data rate for the connection of gamers to their BS through the GaaS slice is set to 10 Mb/s and 50 Mb/s, respectively.

We consider four different scenarios, that correspond to the available options for the association of the two gamers that we monitor in detail to the two edge servers. Scenario (i, j) assumes that g_1 is connected to e_i while g_2 is connected to e_j , with $i, j \in \{1, 2\}$.

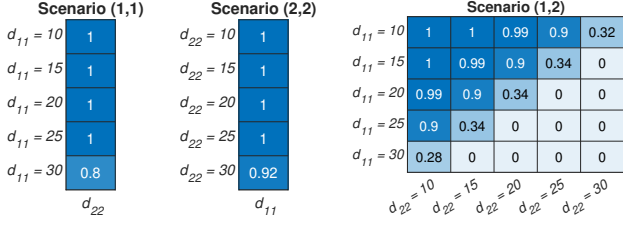


Fig. 3. Success probability with two gamers playing together when VNFs are modeled as M/M/n queues and $\rho=0.9$. Delay d_{ij} in ms.

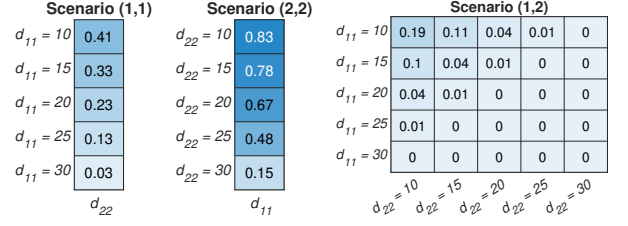


Fig. 4. Success probability with two gamers playing together when VNFs are modeled as M/M/n queues and $\rho=0.99$. Delay d_{ij} in ms.

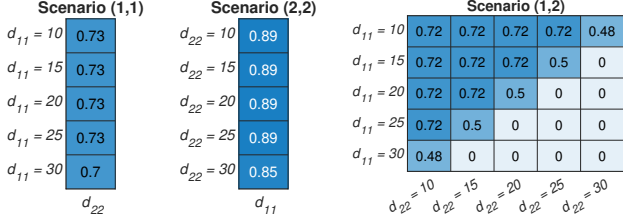


Fig. 5. Success probability with two gamers playing together when VNFs are modeled as M/M/n/k queues and $\rho=0.9$. Delay d_{ij} in ms.

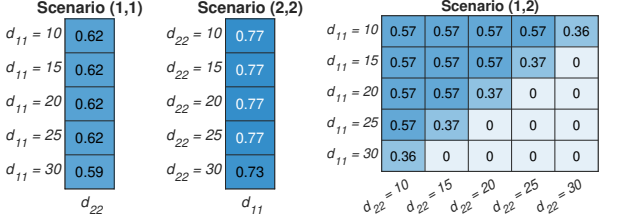


Fig. 6. Success probability with two gamers playing together when VNFs are modeled as M/M/n/k queues and $\rho=0.99$. Delay d_{ij} in ms.

A. Impact of User-Edge Latency

We discuss the impact of latency between gamers and VNFs by reporting in Figs. 3 to 6 success probabilities considering five possible average values for one-way delays, assuming symmetric paths, hence $U_{ij}=D_{ij}$: 10, 15, 20, 25, and 30 ms, in different conditions of edge server load and buffering. On top of the average delay we consider an additive random value taken from a triangular distribution between $t_m=-2$ ms and $t_M=2$ ms, to account for the variable component of the delay incurred in the RAN (whereas the average RAN delay is already counted in U_{ij} and D_{ij}). The timeout is set to $T_O=75$ ms, while the delay between the two edge servers is 10 ms. For a fair comparison of different scenarios, we impose the load of used edge servers to be fixed to either 0.9 or 0.99, which, once we have calculated the traffic due to the two reference gamers g_1 and g_2 , is obtained by appropriately tuning the quantity of background traffic corresponding to other users. Notice that we use as definition of the edge server load the ratio $\rho=\frac{\lambda}{n\mu}$, so that as ρ approaches 1, available computing resources saturate.

Figs. 3 and 4 report the success probability for the two gamers with very large (actually, infinite) edge server buffers. Figs. 5 and 6 show similar results, with $k=2n$, i.e., 5 processors and 5 positions in the buffer at e_1 , and 10 processors and 10 positions in the buffer at e_2 .

Note that in Scenario (1,1) results are invariant with respect to d_{22} , since e_2 is not used; hence, we only report results for variable d_{11} . Conversely, in Scenario (2,2) results are invariant with respect to d_{11} , because e_1 is not used; hence, we only report results vs. d_{22} . In Scenario (2,1) results are invariant to both d_{11} and d_{22} because g_1 connects to e_2 and g_2 to e_1 . Since we set $d_{12}=d_{21}=20$ ms, results for Scenario (2,1) correspond to those of Scenario (1,2) with $d_{11}=d_{22}=20$ ms.

We clearly see that performance degrades with higher latency (i.e. higher distance) between gamers and VNFs, and when approaching the edge server capacity. Both effects are

quite intuitive. We can also observe that very large buffers are useful in some cases, not in others. E.g., Scenario (2, 1) at 90% load performs consistently better with shorter buffers, and almost all scenarios at 99% load perform better with shorter buffers, except for Scenario (2, 2), with $d_{11}=10$ ms.

These results indicate that choosing a buffer size at game VNFs is an important issue. Serving all incoming requests is not useful when approaching saturation because this risks increasing latency for all requests to the point that they miss their deadline. It would be better to instead drop some requests to ensure that those that are served meet their deadline. Not excessively large buffers might be therefore preferable.

B. Impact of Buffers at the Edge

The next set of results looks more closely at the effect of buffering in game VNFs. Figs. 7 and 8 plot the success probability of g_1 and g_2 versus the load of the selected edge servers. The distance between gamers and game VNFs is 20 ms for each gamer-VNF pair, while other parameters are like before. Each figure contains five curves that correspond to (i) no buffer, (ii) a number of positions in the buffer equal to the number of servers ($k=2n$, with $n=5$ at e_1 and $n=10$ at e_2), or (iii) 4 times the number of servers ($k=5n$), and infinite buffer with either (iv) exponential service times, or (v) deterministic service times—with one high-speed processor that, for fairness of comparison, serves requests n times faster than the single processor used with the other types of queues in which n processors are used—and a PS scheduler.

In general, we see that with infinite buffers the success probability drops to zero when $\rho=1$ because queues reach saturation and delays diverge. On the contrary, the presence of a finite buffer yields success probabilities that in some cases remain at acceptable levels, even in (light) overload. This is because the finite buffer generates losses that keep the workload at a feasible level, thus allowing a fraction of the gamers' interactions to reach their destination within the

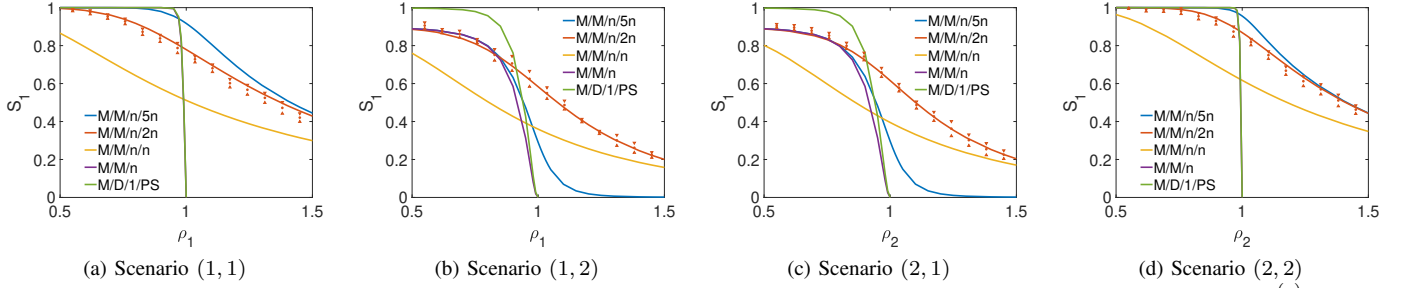


Fig. 7. Success probability of gamer g_1 vs load on the associated edge with $T_O=75$ ms and delays: $d_{ij}=20$ ms, $\forall i, j \in \{1, 2\}$ and $d_s=10$ ms, $d_{12}^{(s)}=20$ ms. Experimental measures with confidence intervals are reported for the M/M/n/2n case

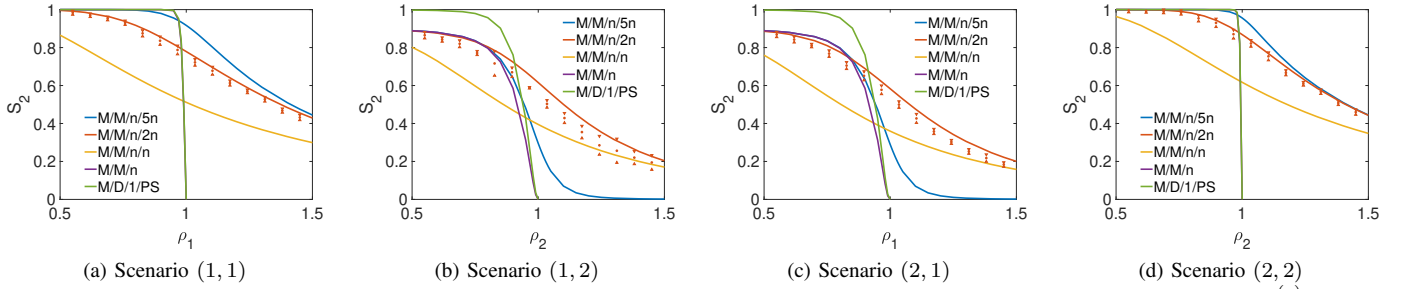


Fig. 8. Success probability of gamer g_2 vs load on the associated edge with $T_O=75$ ms and delays: $d_{ij}=20$ ms, $\forall i, j \in \{1, 2\}$ and $d_s=10$ ms, $d_{12}^{(s)}=20$ ms. Experimental measures with confidence intervals are reported for the M/M/n/2n case

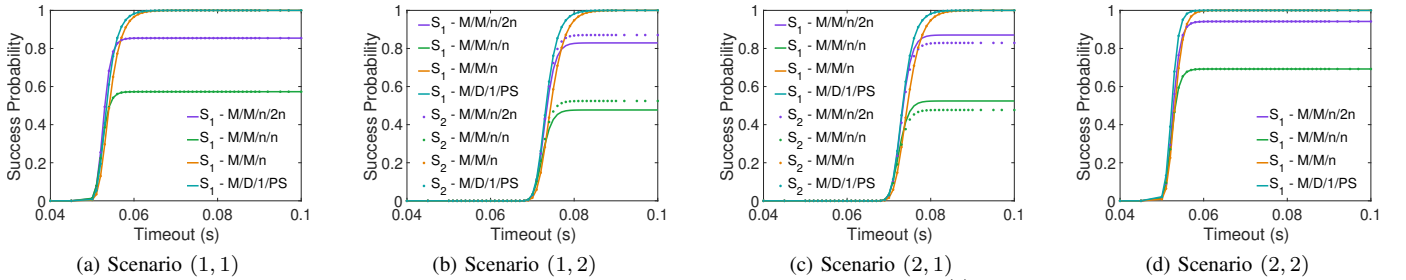


Fig. 9. Success probability of gamers g_1 and g_2 vs timeout, for variable buffer sizes; $d_{ij}=20$, $\forall i, j \in \{1, 2\}$, $d_{12}^{(s)}=20$, $d_s=10$ and load on edges: $\rho_1=\rho_2=0.9$. Results for S_2 in scenarios (1, 1) and (2, 2) overlap those for S_1 , therefore the legend for S_2 is omitted.

specified timeout. For example, in Scenario (2, 2) at load 1.25 both gamers see success probabilities of the order of 0.6. This implies that the gaming service can survive periods of temporary overload with reduced performance, avoiding blackouts. Indeed, slow variations of the number of gamers or of their activity, hence of the VNF load, can translate into the operating point moving along the curves in Figs. 7 and 8.

It is also interesting to observe that, with the chosen parameters, the allocation of both gamers to the same game VNF yields success probabilities equal to ~ 1 for low loads, which is not possible with the allocations to different game VNFs (except for the case of the M/D/1-PS model), due to the additional 20 ms necessary to go from one edge server to the other. The reason why the M/D/1-PS model behaves differently is twofold. First of all, in low load the probability of a small number of customers in the queue is high so that the whole processing capacity is divided among few, each receiving a higher share than with a predefined division in n shares. Second, the deterministic service time avoids long service instances that result from sampling the tail of the exponential distribution. This is on the one hand an indication

of the fact that using fewer powerful processors yields better performance than using many slower ones, and on the other of the fact that assuming exponential service times leads to pessimistic estimates of the overall system performance.

C. Impact of Timeout

The delay with which a gamer visualizes the effect of the actions of gamers that are playing in her game environment is a key QoS metric for online gaming. Curves in Fig. 9 show the success probability of the two gamers in the four possible scenarios with different values for the buffer sizes at the two edge servers (buffer sizes are however proportional to the number of processors at each server). The load of each server is 0.9. Delays from gamers to game VNFs and between game VNFs are like in the previous case.

In Scenario (1, 2) as well in (2, 1) we clearly see that without buffering only about 50% of gamer's commands at most are represented in the video stream within the chosen timeout. Very large (infinite) buffers allow success probabilities to reach one, but only for timeouts of the order of about 85 ms or more (with little benefit deriving from deterministic service times).

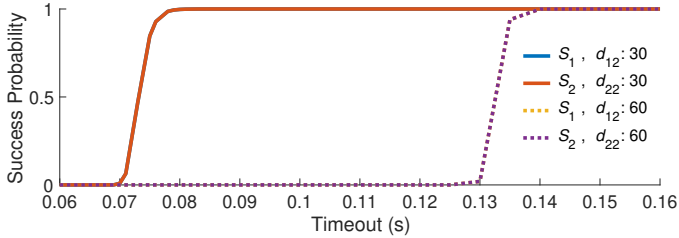


Fig. 10. Cloud gaming scenario versus edge gaming scenario with server load $\rho=90\%$, $d_s=10$ ms; servers modeled as $M/M/n$ queues with $n=50$ in the cloud and $n=10$ in the edge.

In the two scenarios in which a single game VNF is used, timeouts are easily met in the very large buffer cases for values as low as 60 ms. On the contrary, with small or no buffer, the loss at edge servers prevents success probabilities to reach desirable levels, with the possible exception of Scenario (2, 2), provided a success probability around 95% is acceptable.

VI. CLOUD VERSUS EDGE GAMING

Our final set of results is reported in Fig. 10. Curves refer to an edge gaming scenario in which both gamers are connected to a GaaS slice hosting a game VNF on the edge (an “edge VNF”) and a cloud gaming scenario, where both gamers connect to the same “cloud VNF”. The distance of both gamers from the edge VNF is 30 ms, and 60 ms from the cloud VNF. The load of servers is 90%. The number of processors at the edge VNF is 10 whereas at the cloud the number is 50. Buffers are infinite in both cases. Note that fixing the load at the two servers implies that the edge performs 20% of the work done by the cloud, but this is reasonable since an edge architecture replaces few (very powerful) cloud servers with many (less powerful) edge servers. This implies a more distributed load over edge servers for a given population of gamers.

We clearly see that the edge VNF cannot satisfy timeout values less than 71 ms (30+30 to go and come back to the game server, 10 for video preparation, 1 for processing at the edge), while the cloud VNF cannot provide latency lower than 131 ms (60+60+10+1). Success probabilities grow rapidly beyond those latency values and reach close to 1 at around 80 ms in the case of the edge VNF, and 140 ms in the cloud case. Note also that the symmetry of the scenario is such that curves overlap for the two gamers in both edge and cloud cases.

The advantage of the edge architecture is intrinsic in the distribution of computing facilities over the network, whose objective is to reduce the latency between users and computing facilities, and is especially beneficial for latency-constrained services, like the gaming application we are considering.

VII. VALIDATION

A. Measuring GaaS in the wild

To assess the GaaS performance in real-world edge computing scenarios, we set up a distributed architecture akin to the one described in Fig. 1. Specifically, we developed a flexible set of instrumented microservices providing the building blocks of our architecture. Three main blocks were defined: *client*, *server*, and *UPF*. The setup architecture is depicted

in Fig. 11. The micro-services comprise QUIC endpoints, are coded in Golang, and are containerized using Docker.

The *client* container consists of several elements contributing to setting up an active connection with the game server and also to tracing and measuring packet receptions, losses, and time statistics. The packet generator is responsible for modulating the traffic flow, which subsequently gets marked with the specific type—either control (*Ctrl*) or feedback (*Fb*), respectively. Moreover, for the *Ctrl* flow, the marker adds a timestamp, the client identifier, and a unique identifier for the packet. Then, dedicated *goroutines* (i.e., lightweight Golang threads) handle packet transmissions through QUIC connections; these *goroutines* are also responsible for collecting data and computing statistics from the incoming packets by matching identifiers, with the help of a list of *pending* control threads (i.e., control packets already sent, whose incorporation in the game stream is not yet completed). Uplink traffic passes through the *UPF* container, emulating the behavior of a tunnel between the player and the game server with a QUIC connection. Specifically, the *UPF* implements a single server with a finite but huge buffer size. At this point, requests are routed toward the game server associated with the player originating the traffic flow, obtaining the network setting of Fig. 11. Packets incoming to the *Edge* containers are either local—that is, incoming from the associated player—or *Inter-edge control* packets coming from the other game server, see the solid and dashed lines in Fig. 11 between the two edge servers. Local packets get to the QUIC endpoint where *Fb* packets get discarded, whereas *Ctrl* packets originate new *Inter-edge control* packets (that are necessary for the remote VNF to generate the video to be delivered to the remote gamer). Also, identifiers and timestamps associated with *Ctrl* packets get inserted in packets carrying the video stream, marking the fact that such frames report the updates related to players’ control instructions, as depicted at the bottom of Fig. 11. Similarly, packets incoming from *Inter-edge control* flow get integrated into the downlink stream, with the only difference being that they are associated with a different QUIC endpoint. These timestamps and identifiers are necessary for the computation of the delays with which game commands are incorporated in the video stream observed by players.

B. Experimental results

We validated our model by comparing model predictions and experiment measurements with two game VNFs positioned at a distance of about 20 ms. Both gamers experienced similar latency when connecting to their nearest game VNF, also situated about 20 ms away.

Specifically, the game VNFs are 19.30 ms away from each other, with a standard deviation (σ) of 2.23 ms. The two players are 20.02 ms ($\sigma=1.76$ ms) and 19.77 ms ($\sigma=2.55$ ms) away from their associated game VNF, respectively. The time required to implement changes to the video stream in the game VNF is set to $d_s=10$ ms. Game VNFs have either 5 or 10 processors and the same number of positions in the buffer. Gamers access their game VNF through their associated

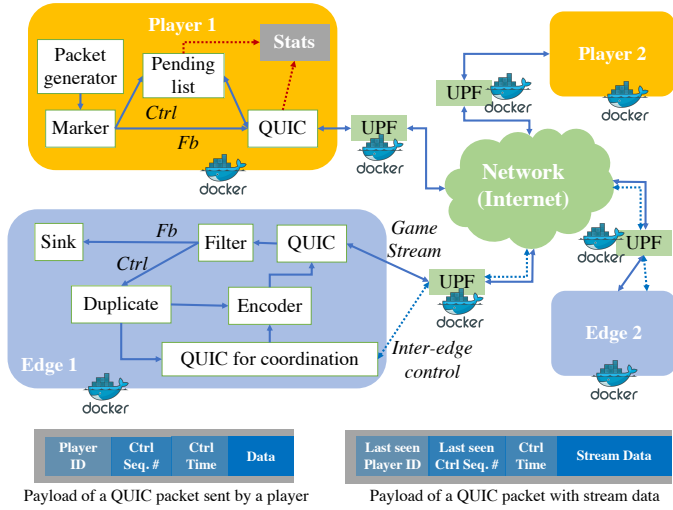


Fig. 11. Experimental setup. To enable the computation of statistics, control (Ctrl) and stream packets carry the relevant ID of players and control sequence numbers, in addition to timestamps.

UPF, requiring constant times equal to 0.1 ms to process uplink packets and 0.2 ms for the downlink ones. Gamers generate an uplink stream of 1 000 control packets per second of 1 000 bytes each, with only one in every one hundred packets delivering commands to the game streamer.

We evaluate each gamer’s success probability considering a latency threshold of 75 ms with different background traffic intensity; in our experiments, the background traffic ranges from $\rho_i=0.5$ to 1.5, with a step size of 0.1. The comparative evaluation of our model with real-world measurements is illustrated in Figs. 7 and 8, accompanied by 90% confidence intervals. We can observe a very good match between model predictions and experiments in spite of the many simplifications introduced in the model.

VIII. CONCLUSIONS

We have proposed an online gaming service deployed on edge computing facilities and studied the performance of the system. Through analysis and experiments, we have highlighted how our proposal is, on the one hand, very promising, while being, on the other hand, also very challenging. Indeed, the downside of the edge option is the need to coordinate gaming engines on different edge facilities and VNFs.

Our analytical results show that, for a number of parameter values that reflect real gaming systems, the edge gaming option can yield lower latency than cloud gaming, provided that the allocation of interacting gamers to edge servers is carefully chosen. It often happens that better performance is achieved by allocating gamers that closely interact on the same edge server, but under some combinations of the system parameters, an allocation of gamers to their closest game VNF can be more effective, in spite of the need for interaction among servers.

Our model can provide a very effective tool for managers of GaaS applications since it enables the choice of the best performance option upon the arrival of new groups of gamers. This allows the GaaS slice to operate at its best performance and guarantee the best possible QoE to gamers.

In the future, we plan to analyze more realistic gaming workloads as well as more realistic networking features, with the constraints imposed by connect-compute architectures that already exist or are under study in the 3GPP galaxy.

REFERENCES

- [1] European Commission, “European Media Industry Outlook report,” <https://digital-strategy.ec.europa.eu/en/library/european-media-industry-outlook>, 2023.
- [2] ITU. SG12: Performance, QoS and QoE.
- [3] S. Möller *et al.*, “Gaming taxonomy: An overview of concepts and evaluation methods for computer gaming QoE,” in *QoMEX*, 2013.
- [4] ITU, “ITU-T Recommendation ITU-T G.1032, Influence factors on gaming quality of experience,” ITU, Tech. Rep., 2017.
- [5] —, “ITU-T Recommendation ITU-T P.809, Subjective evaluation methods for gaming quality,” ITU, Tech. Rep., 2018.
- [6] —, “Recommendation ITU-T G.1072, Opinion model predicting gaming quality of experience for cloud gaming services,” ITU, Tech. Rep., 2020.
- [7] K.-T. Chen *et al.*, “On the Quality of Service of Cloud Gaming Systems,” *IEEE Trans. on Multimedia*, vol. 16, 2014.
- [8] A. Di Domenico *et al.*, “A Network Analysis on Cloud Gaming: Stadia, GeForce Now and PSNow,” *Network*, vol. 1, 2021.
- [9] O. S. Peñaherrera-Pulla *et al.*, “Measuring Key Quality Indicators in Cloud Gaming: Framework and Assessment over Wireless Networks,” *Sensors*, vol. 21, 2021.
- [10] M. Carrascosa *et al.*, “Cloud-gaming: Analysis of Google Stadia traffic,” *Computer Communications*, vol. 188, 2022.
- [11] X. Marchal *et al.*, “An Analysis of Cloud Gaming Platforms Behaviour Under Synthetic Network Constraints and Real Cellular Networks Conditions,” *J. Netw. Syst. Manage.*, vol. 31, 2023.
- [12] Y. Kim *et al.*, “E-Render: Enabling UHD-Quality Cloud Gaming Through Edge Rendering,” *IEEE Access*, vol. 10, 2022.
- [13] G. Z. Papadopoulos *et al.*, “IETF Reliable and Available Wireless (RAW): Use Cases and Problem Statement,” in *Ad-Hoc, Mobile, and Wireless Networks*, L. A. Grieco *et al.*, Eds. Springer
- [14] M. Bacco *et al.*, “A Simulation Framework for QoE-Aware Real-Time Video Streaming in Multipath Scenarios,” in *Ad-Hoc, Mobile, and Wireless Networks*, L. A. Grieco *et al.*, Eds. Springer
- [15] L. De Giovanni *et al.*, “Revamping Cloud Gaming With Distributed Engines,” *Internet Computing*, vol. 26, 2022.
- [16] X. Zhang *et al.*, “Improving Cloud Gaming Experience through Mobile Edge Computing,” *IEEE Wireless Comm.*, vol. 26, 2019.
- [17] F. Spinelli *et al.*, “Edge Gaming: A Greening Perspective,” *Comp. Comm.*, vol. 192, 2022.
- [18] M. Ghobaei-Arani and et al, “An autonomous resource provisioning framework for massively multiplayer online games in cloud environment,” *J. of Netw. and Comp. Appl.*, vol. 142, 2019.
- [19] Y. Gao *et al.*, “Energy- and Quality of Experience-Aware Dynamic Resource Allocation for Massively Multiplayer Online Games in Heterogeneous Cloud Computing Systems,” *IEEE Trans. on Services Computing*, vol. 16, 2023.
- [20] Y. Liang *et al.*, “Interaction-Oriented Service Entity Placement in Edge Computing,” *IEEE Trans. on Mobile Computing*, vol. 20, 2021.
- [21] Z. Chen *et al.*, “Wireless Multiplayer Interactive Virtual Reality Game Systems With Edge Computing: Modeling and Optimization,” *IEEE Trans. on Wireless Communications*, vol. 21, 2022.
- [22] S. Agarwal *et al.*, “VNF Placement and Resource Allocation for the Support of Vertical Services in 5G Networks,” *IEEE/ACM Trans. Netw.*, vol. 27, 2019.
- [23] G. Einziger *et al.*, “Virtual Service Embedding with Time-Varying Load and Provable Guarantees,” *IEEE Trans. on Cloud Computing*, vol. 11, 2022.
- [24] E. T. S. Institute, “TS 123 501 - V15.3.0 - 5G; System Architecture for the 5G System (3GPP TS 23.501 version 15.3.0 Release 15),” ETSI, Tech. Rep., 2018.
- [25] T. J. Ott, “The Sojourn-Time Distribution in the M/G/1 Queue with Processor Sharing,” *J. of Appl. Prob.*, vol. 21, 1984.
- [26] Y. Pan *et al.*, “The First 5G-LTE Comparative Study in Extreme Mobility,” *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 6, 2022.
- [27] R. A. K. Fezeu *et al.*, “An In-Depth Measurement Analysis of 5G mmWave PHY Latency and Its Impact on End-to-End Delay,” in *PAM Conference, 2023*, 2023.