

Received August 8, 2021, accepted August 16, 2021, date of publication August 24, 2021, date of current version September 3, 2021.

Digital Object Identifier 10.1109/ACCESS.2021.3107624

# QoS-Based Data Aggregation and Resource Allocation Algorithm for Machine Type Communication Devices in Next-Generation Networks

NASSER AHMED<sup>1</sup> AND NASSER-EDDINE RIKLI<sup>1</sup>, (Senior Member, IEEE)

Department of Computer Engineering, College of Computer Science and Information System (CCIS), King Saud University, Riyadh 11451, Saudi Arabia

Corresponding author: Nasser Ahmed (nahmed@ksu.edu.sa)

This work was supported by the Grant of the Research Center of the College of Computer and Information Sciences, Deanship of Scientific Research, King Saud University.

**ABSTRACT** Machine Type Communication (MTC) becomes one of enablers of the internet of things, it faces many challenges in its integration with human-to-human (H2H) communication methods. To this aim, the Long Term Evolution (LTE) needs some adaptation in the scheduling algorithms that assign resources efficiently to both MTC devices (MTCDs) and H2H users. The minimum amount of LTE resources that can be assigned to one user is much larger than the requirements of a single MTCD. In this paper, a QoS-enabled algorithm is proposed to aggregate MTCD traffic coming from many sources at the Relay Node (RN) that classifies and aggregates the MTCD traffic based on the source type and delay requirements. In this study, three types of MTCD and one H2H sources will be considered. Each type of MTCD traffic will be grouped into a separate queue, and will be served with the appropriate priority. Resources are then assigned to the aggregated MTC traffic instead of an individual assignment for each MTCD, while the H2H users will be directly connected to the LTE. Two schemes of resource partitioning and sharing between the MTCDs and the H2H users will be considered: one proportional and the other moving-boundary. Simulation models will be built to evaluate the proposed algorithms. While the obtained results for the first scheme showed a clear improvement in LTE resource utilization for the MTCDs, a negative effect was noticed in the performance of the H2H users. The second scheme achieved a positive improvement for both MTCDs and H2H users.

**INDEX TERMS** 5G networks, data aggregation, Internet of Things (IoT), MTC, resource allocation, quality of service (QoS).

## I. INTRODUCTION

An increasing demand for high data rates, high capacity, and low latency to support a fully connected networked society that offers access to information and the sharing of data anywhere and anytime for anyone and anything, has led to the introduction of a new type of communication paradigm called machine-to-machine communication (M2M) or machine type communication (MTC). This type of communication implies that machines have the ability to communicate with each other in a smart approach without or with a minimum of human intervention [1]. Interest in MTCDs has increased in recent decades because they exist

The associate editor coordinating the review of this manuscript and approving it for publication was Miguel López-Benítez<sup>1</sup>.

in many applications of the internet of things, such as but not limited to e-Healthcare, smart metering, smart cities, intelligent transportation systems, supply chains, surveillance monitoring systems, the prediction of natural disasters, and many social applications [2].

LTE-Advanced (LTE-A) is a candidate as the most suitable cellular technology to support MTC, due to its high data rates, large coverage area, high capacity, and spectrum efficiency. However, there are many challenges in the integration of MTCDs in an LTE-A network [3]–[5]. LTE-A has largely been designed to support H2H devices, which typically require high data rates and a small delay, have a small number of users (compared to MTCDs) and transmit a large volume of data packets. In contrast, MTCDs have different characteristics, such as a large number of devices, a low data

rate, a small data packet size, upload-centric applications, and power constraints [5]. This contradiction between the characteristics and requirements of H2H and MTCDS is considered one of the biggest challenges in the use of LTE-A. Cisco estimates that the number of MTCDS globally will increase to 14.7 billion by 2023, and they represent 50% of all connected devices [6]. A large number of MTCDS trying to get access simultaneously to a base station (BS) in an LTE-A system brings another challenge to the integration of MTCDS into LTE-A networks.

Radio resource allocation is one of the largest challenges facing the integration of MTCDS in an LTE-A system. The main difficulty is inefficient resource allocation, which is due to several factors. Firstly, H2H and MTCDS have quite different characteristics. H2H traffic is download-dominant, with a small number of users and a large data packet size. MTCDS, on the other hand, have upload-dominant traffic and a huge number of devices with a small data packet size. Second, in LTE-A, the minimum amount of resource blocks (RBs) that can be allocated to one User Equipment (UE) exceeds the requirements of MTCDS. For example, the smallest amount of RBs that can be allocated to one UE in LTE/LTE-A is one physical resource block (PRB), which contains  $(12 \times 7)$  resource elements. This can be used to transmit hundreds of bits of data; however, the requirement of most MTCDS does not require this amount of resource due to the small size of their packets. This makes it inefficient to assign one PRB to one MTCDS. Therefore, a new mechanism should be designed to manage radio resource allocation for MTCDS in LTE-A systems in a more efficient manner, without creating negative effects for H2H traffic. The third challenge is power consumption due to the power constraints in MTCDS, in particular when it is difficult to recharge the battery of the MTCDS, or it cannot be recharged, when the MTCDS are placed in a critical environment. Therefore, MTCDS require efficient power management.

Data aggregation is one of the most practical solutions used to solve the problem of resource allocation to MTCDS. This is achieved by clustering and multiplexing MTCDS traffic from many MTCDS into an aggregator, which in turn sends the aggregated data to the next stage. This aggregator has a powerful capability in terms of energy, computation, and storage; it may be a cluster head of a capillary network, or it may be a cellular-based design within an LTE RN. The issues of aggregation, multiplexing, and resource allocation have been examined by many researchers, as can be seen in the related works in Section II.

In this paper, a QoS-based data aggregation algorithm is proposed for LTE-A networks, which incorporate both MTCDS and H2H users. The proposed data aggregator is cellular-based, and it has been designed within the LTE-A RN. It aggregates data from different types of MTCDS with different types of QoS requirements then classifies the traffic based on their QoS to different queues, buffering the aggregated traffic until an adaptive time threshold or until an adaptive buffer size. At that point, the aggregator implements

a frame formulation and multiplexing technique by accumulating the traffic from each buffer, based on their priority, into a new, large LTE frame, and then transfers the accumulated large frame to LTE Evolved Node Base Station (eNB). Therefore, the LTE eNB assigns resources to the aggregator RN instead of individual MTCDS.

The rest of this paper is organized as follows: the next section presents the related works, and the contribution of this paper is introduced in section III; Section IV introduces the system model, resource management in the proposed schemes is presented in section V; section VI presents the aggregation function; the performance metrics are defined in Section VII; the simulation model and configuration are presented in Section VIII; the analysis of the results is introduced in Section IX; and, finally, references are listed.

## II. RELATED WORKS

The LTE resource allocation has been studied extensively in recent years, and many approaches have been proposed in the literature. They can be divided into two main categories based on using or not using aggregation for MTCDS in the resource management process. The next subsection explores the first category solution, while the second subsection explores the second category, and finally the third subsection presents the data aggregation based on Software Defined Networks (SDN) and Fog Computing.

### A. RADIO RESOURCE ALLOCATION FOR M2M WITHIN LTE-A WITHOUT DATA AGGREGATION

In this category, the proposed solutions include relaying MTCDS to the eNB while connecting H2H users directly to the eNB, with an orthogonal resource partitioning between the access link and backhaul link [7], an Energy-aware radio resource management (RRM) [8], an energy-efficient resource allocation algorithm with the objective of maximizing bits-per-joule [10]. A context-aware resource management approach for MTCDS gateways is proposed in [9], while in [10], a delay-aware radio resource scheduler algorithm that satisfies the QoS requirements for MTCDS and H2H is presented, and a hierarchical RRM approach is proposed in [11]. In [12], a type-2 fuzzy logic controller mechanism is used for radio resource allocation for MTCDS in co-existence with H2H within LTE, and a real-time spectrum analyzer is used for resource management in [13]. A tree-based algorithm is used in [14], and a Maximum energy efficiency is investigated in [15]. Each of these studies is explored in more detail below.

In [7], the authors propose a radio resource partition pattern for the downlink transmission of LTE-A cellular networks with MTCDS communications. Multi-hop transmission is defined for MTCDS, which are connected through a machine type communication gateway (MTCG) to the eNB to mitigate the massive competition for radio resources. MTCDS to MTCG and MTCG to eNB links are assigned orthogonal parts of the radio resources, while all other links are directly associated with the eNB and share the remaining resources of the channel. A user utility function was defined in terms of

the achievable data rate, and though its maximization the corresponding radio resource allocation matrix was determined.

The limitations of [7] in the cases of low traffic rate and delay tolerant features of MTC, were addressed in [8] by presenting an energy-aware RRM scheme for MTC/H2H co-existence scenarios in LTE networks, with guaranteed QoS requirements for different users. This was achieved through minimization of the overall transmission power and maximization of tolerable packet delay for MTC. Two heuristic algorithms based on the steepest descent approach were proposed to solve this optimization problem. The first shows how to effectively achieve the goal of transmitting H2H and MTC data at the minimum power, while the second takes into account only the minimization of transmission power for H2H traffic.

The authors in [16] extend the work in [7] further by proposing an energy-efficient resource allocation algorithm with the objective of maximizing bits-per-joule capacity under statistical QoS provisioning. The proposed scheme was analyzed using mixed-integer programming, and the optimization problem was solved with canonical duality theory.

A context-aware resource management approach for MTC gateways was proposed in [9] to achieve QoS provisioning by analyzing data on the traffic flow generated by H2H and MTC users. Various classes of H2H/MTC traffic were considered, namely: conventional, streaming, interactive, background, priority alarm, time tolerant, and time controlled. Also, dynamic contextual information was taken into consideration, such as service type, MTC type, and network status, and then the MTC services were adapted to these diverse contexts. The main achievements were a mitigation of congestion and overload conditions in the system by satisfying the MTC services without degrading QoS for existing H2H services.

In [10], the authors proposed a delay-aware radio resource scheduler algorithm, which satisfies the QoS requirements for MTC while ensuring a minimal impact on the QoS of H2H traffic. The MTC and H2H flows are grouped into  $n$  different classes according to their remaining times to serve (RTTS), defined as the time within which the flow should be served by the scheduler to meet its delay tolerant time. The RBs are assigned to classes according to a priority that is inversely proportional to the RTTS values. Moreover, within the same class, the scheduler gives a higher priority to H2H over MTC to avoid the negative impact of MTC on H2H. Although this approach satisfies the QoS requirements of each flow in terms of delay and data rate, the grouping of MTC and H2H devices is managed at the traffic flow level. There is no grouping for the device itself, no details about the location, the mobility, or the power consumption. In addition, this approach assumes direct access between MTC and the eNB, which is not suitable for a massive number of devices. Moreover, starvation may occur for the delay tolerable MTCs in the case of high congestion.

In [11], the authors propose a hierarchical RRM approach. As in typical MTC applications, the amount of data consumed is relatively small, the RBs granted to MTC

are not fully consumed. Consequently, C-UEs can exploit this unused portion that would otherwise be wasted. In the proposed scheme, a two level hierarchy is proposed. In the first level, a PRB is allocated to MTC as well as to C-UE, while in the second, the MTC delegates a portion of its unused resources to a neighboring C-UE. The results showed that, in the case of the high load of MTC, limited gain was achieved.

In [12], the authors present a radio resource allocation mechanism in LTE for MTCs co-existing with H2H devices and using a type 2 fuzzy logic controller. They assume an ideal channel where the failure of any access request can only occur as a result of its collision. Two categories of applications were considered: real-time (RT) applications, which are sensitive to delay, and non-real-time (NRT) applications, which are delay tolerant but have a minimum power requirement. This mechanism consists of two stages. In the first stage, the system evaluates the data flow based on the decision factors, while in the second step, RBs are allocated by first assigning them to RT users and then assigning the remaining RBs to NRT users.

The impact of different channel conditions on radio resource utilization in real LTE networks was analyzed in [13]. A commercial RT spectrum analyzer was used to analyze the uplink LTE resource utilization, which was computed as a function of the number of RBs, as well as the data rate and spectrum efficiency. The main goal was to minimize the impact of MTC traffic on H2H traffic, which were co-existing on the same LTE network. This was achieved by allowing the MTCs to transmit data on the channel with both high probability and high quality.

Another variant using a persistent resource allocation algorithm for MTC was proposed in [14]. The resources of the MTCs were allocated periodically in a recursive manner based on a tree structure. This scheme does not use any resource for RACHs; instead, it assigns all resources as uplink data channels without any additional control signaling during the life of a machine. The concept of the persistent resource allocation scheme was to multiplex as many machines of different periods as possible onto a single channel. The tree-based algorithm was used to determine if the state of machines with different periods can be multiplexed. This scheme has shown potential performance gains in supporting a larger number of devices in comparison to coordinated access schemes for small packet transmissions. However, it was only beneficial for periodic traffic and was not useful for aperiodic or bursty data.

In [15], the authors investigated the maximum energy efficiency of MTC data packet transmission with the uplink SC-FDMA in LTE-A. They formulated the problem of energy efficiency as an optimization problem that includes modulation and coding scheme assignment, resource allocation, power control, and other constraints in the uplink of an LTE-A network. The problem was then converted into an NP-hard mixed-integer linear fractional programming problem to reduce the computation complexity and find the final

optimum level of energy efficiency. They assumed different types of MTCD with different types of sensors generating different types of data packets. In this way, it was not possible to aggregate data into one large packet, since each sensor has to report its data in a determined time interval. The results of the simulation showed that, with limited RBs, the proposed algorithm achieved a low packet dropping rate with optimal energy efficiency in the case of large number of MTCDs.

## B. RESOURCE ALLOCATION FOR M2M DEVICES USING DATA AGGREGATION

The second category of research into resource allocation and management for MTCDs in a co-existent network with H2H users, explores the research undertaken in data aggregation and multiplexing for MTCD. Data aggregation can be achieved in three ways:

- 1) Data aggregation at the MTCD level, in which the MTCD delays and aggregates its data by itself before transmitting it to the eNB. This method can be used to increase the efficiency of resource allocation. However, in most cases, MTCD traffic flows are periodic, sending their data at predetermined intervals and with only a small amount of data sent at each time interval, which makes this solution impractical.
- 2) The regular H2H mobile users can be used as mobile aggregators to aggregate and attach the MTCD traffic to its own data using its own unused resources [17], [18]. This can increase resource utilization by exploiting the unused resources of the traditional user, which otherwise would be lost. However, this solution is not suitable for high priority MTCD traffic that cannot wait for the availability of unused resources assigned to traditional users.
- 3) The most practical solution is the aggregation, clustering, and multiplexing of MTCD traffic from many devices into an aggregator (cluster head/gateway/RN), which in turn transmits the aggregated traffic to the LTE eNB which assigns its resources to the aggregator node instead of individual MTCD. This solution needs a number of algorithms to manage resource allocation, address how to aggregate the MTCD flows into one node, handle multiplexing issues, manage power consumption, and select the appropriate aggregator.

The benefits of data aggregation are not only in resource allocation efficiency, but also in other areas such as reducing power consumption [19], [20], increasing system capacity, increasing the scalability of the system to serve a massive number of MTCDs, and decreasing the signaling overhead [21], [22]. Much research has been conducted in relation to data aggregation, clustering, and multiplexing [23]. Data aggregation can be categorized in terms of the type of aggregator as either fixed data aggregator (FDA), or mobile data aggregator (MDA), or cooperative data aggregation (CDA).

Alternatively, data aggregation can be classified based on radio access technologies into two types: cellular-based aggregators and capillary-based aggregators. In the former,

the MTCDs are equipped with a subscript identity module and connected to the network through the cellular gateway using a licensed frequency band [16], [24], [25]. In the latter, MTCDs are connected to the network through a capillary gateway using an unlicensed frequency band (e.g., ZigBee or Bluetooth Low Energy), while the aggregator itself is connected to the BS using a licensed band such as LTE-A [22], [23], [26], [49]. As this classification has been the mostly accepted and used, we will present the two categories in more details.

### 1) FIXED DATA AGGREGATOR

Fixed data aggregators (FDA) can be further categorized into two types: single fixed data aggregator (SFDA), and multiple fixed data aggregators (MFDAs). In the former type, only one aggregator is used, while in the latter, many aggregators are used. In a single data aggregator, the signaling overhead between the aggregator and the eNB is reduced; however, the risk of being a single point of failure is increased. In addition, single data aggregators increase the delay of aggregated packets, and MTCDs may overwhelm the aggregator with huge numbers of packets, increasing the ratio of dropped packets. In contrast, using multiple RN aggregators increases the signaling overhead between MTCDs and the eNB, but it provides more reliability.

Single data aggregators were proposed in [27]–[29]. In [27], the small data packets from MTCDs are aggregated, delayed, multiplexed, and reformatted to a large packet at the Packet Data Convergence Protocol (PDCP) layer within the RN. Resource utilization improved at the cost of delay. In contrast, a hierarchical energy-efficient data aggregation model for MTCD uplink to minimize the average energy density consumed was proposed in [28], where a multi-stage and a hierarchical structure were used to select some MTCDs in a probabilistic way to work as aggregators to the data packets from other nodes. At each stage, there is a new hierarchy of aggregators that receives data from the aggregators of the previous stage. Finally, in [29] a data aggregation for massive MTC in a large-scale cellular network was introduced. The authors investigated the signal to interference ratio (SIR) for both the aggregation phase and the relaying phase. They also analyzed the performance of the system in terms of the average number of successful MTCDs and the probability of successful channel utilization using a stochastic geometry framework. Two resource scheduling approaches were used: a random resource scheduling (RRS) algorithm, and a channel-aware resource scheduling (CRS) algorithm. The results showed that the CRS algorithm outperforms the RRS algorithm.

The MFDAs scheme were presented in [30]. Here, the MTCD can be connected to one or more MTCGs at the same time. Two types of relaying techniques were introduced. In the first, an SIR based, the signal from the MTCD can be decoded by one or more MTCGs; therefore, the packet may be duplicated at the eNB. In the second, a location-based, the packet duplication drawback was overcome by

allowing the MTCs to transmit only to the closest MTCs. This improvement was accomplished at the cost of increasing the information exchanged between the MTC and MTCs. This work was applied only to homogeneous types of MTCs with the same type of traffic, and the QoS and delay tolerant MTC services were not taken into account.

## 2) MOBILE DATA AGGREGATOR

In a Mobile Data Aggregator (MDA), one or more mobile data aggregators were used to first aggregate the data from the MTCs and then relay it to the eNB. The mobile data aggregator can be a mobile RN installed on a mobile vehicle (e.g., public bus, taxi), a UE that allows MTCs to connect and send their data through it [17], [18], or an RN installed on a drone/mobile unmanned aerial vehicle (UAV) [31], [32]. Because of the mobility of MDAs, when they enter the vicinity of MTCs and allow the MTCs to connect and send their data through them, they reduce the communication distance between the MTC and MDA gateway, and thus decreasing the transmission power needed. This scheme is best suited to the aggregation of periodic and delay tolerant MTC traffic [18], such as smart metering, due to the fact that the MTC has to wait for the MDA to arrive at its trajectory during its journey.

The use of a UE as an MDA has been introduced in many research studies in the field [17], [33], [34], although some researchers prefer not to use a UE as an MDA because it causes fast depletion of the UE's battery. Some authors have suggested implementing energy harvesting for mobile UE to overcome battery depletion issues [35]. Multiplexing the bandwidth between MTCs and regular UEs has been proposed by 3GPP Release 13 and beyond, so that the MTC traffic can be trunked and multiplexed within the resources assigned to regular Device to Device (D2D). Using only one gateway as an MDA is referred to as a single mobile data aggregator (SMDA), while using more than one gateway as an MDA is referred to as a multiple mobile data aggregator (MMDA).

Using a UE as SMDA has been proposed in [17], [33], [34], [36]. In [17], the conventional UE is used as a single mobile gateway aggregator, and the communication between D2D is exploited in the cellular system to aggregate and multiplex the traffic from surrounding MTCs. The UE attaches its own data and then uses a Time Division Multiple Access (TDMA) to relay all data to the eNB. Through this method, the mobility of regular D2D is exploited to decrease the transmission distance between MTCs and the eNB, thereby decreasing the power consumption of MTC transmission. It also mitigates the capacity drawback in the large-scale system by grouping the MTCs to regular users. Its drawback however, is the increase in MTC traffic delay. In [33], the authors use two applications to investigate the potential usage of the smartphone as a mobile gateway for MTCs using standard middleware. They show improvement in system connectivity but at the cost of smartphone battery depletion and increased delay for MTC traffic. In [36],

the authors propose a scheme for MTC traffic aggregation and trunking within the resources of D2D users in a large-scale system. They introduce a comprehensive stochastic geometry framework to analyze the coverage area of regular users, to make sure that the MTCs send their data using the shortest path to a nearest regular user. The model assumes that an MTC is connected to only one UE to ensure that the aggregation process is achieved in a distributed manner.

Multiple mobile data aggregators (MMDAs) have been proposed in [31], [32], [41]–[43]. In [31], the authors proposed a resource allocation and scheduling scheme for cluster-based MTC. The goal was to increase the power efficiency of the system while meeting the rate requirement for each MTC device. Each MTC group had a cluster head (CH) that worked as both coordinator and aggregator to collect data packets from the MTCs and send them to a flying BS on a UAV. Orthogonal frequency division multiple access (OFDMA) was used for uplink, and the queue rate stability approach was used to determine the minimum number of required UAVs to serve the CHs. Although this study showed good results in terms of power consumption for the CH and the minimum number of UAVs required, it required other protocols and algorithms such as obtaining the positions of CHs and computing the dwell time of UAVs over the CHs. The work in [32] is an extension to [31], where an efficient deployment and mobility model for the UAVs was introduced. The mobility of UAVs was determined and the power consumption by the UAVs was minimized, while the MTCs were also served with minimum transmission power.

In [37], co-existing H2H users and MTCs were considered with the H2H users acting as MDAs to collect data from the MTCs within their vicinity and relay it to the eNB. The resources for MTCs were allocated based on the residual energy in the MTC: high priority was given to the MTC with less residual energy. Results showed that the delay constraints for both H2H and MTC were satisfied, and an improvement in system performance in terms of energy efficiency was achieved, thereby extending the network lifetime.

In [38], the authors introduce a stochastic geometry-based framework to analyze the coverage probability and average data rate of a three-hop MTC distributed in co-existence with regular UE (H2H users). The UEs were used to relay the data of MTCs in multi-hops to the eNB without aggregating data from different MTCs. The results showed an improvement in terms of data rate and network area coverage, due to the fact that MTCs out of range can be relayed using UE by exploiting D2D links. The mobility of UE was addressed by using a space-time graph to predict the location of UEs and exploit it to design a cost efficient multi-hop D2D topology. Good results were achieved in terms of data rate and extending the coverage area of the network. However, the study did not take into consideration the transmission delay, which should have been taken into consideration.

The work in [17] was extended in [39], with the proposal of three aggregation schemes: one fixed, one random, and another greedy. In all these schemes, the UE is used as an

aggregator gateway to aggregate the traffic from the MTCDs and then relay them to the eNB. The authors introduce a mathematical model to evaluate the end-to-end outage probability for the uplink data at the UEs. They show that the greedy scheme outperforms the other schemes in term of outage probability at the MTCD.

A load balancing relay algorithm is introduced in [40], in which mobile MTCDs are grouped randomly and their data is aggregated to an MTC gateway. The MTCDs are regrouped based on the load of each gateway to balance the load and resources for each gateway. Dynamic resource allocation for MTCDs in the link between MTCD and MTCG is studied and system performance is evaluated in terms of system capacity and outage probability. The results show good performance. However, the authors assume an information exchange (e.g., location information, grouping decision) between MTCDs and BS to achieve the dynamic grouping of MTCDs, where the decision of grouping is assigned by the data aggregation center at the BS, which results in a huge signaling overhead in the backhaul link. Furthermore, QoS is not included in this study.

### 3) COOPERATIVE DATA AGGREGATION (CDA)

While MDA is suitable only for tolerable delay traffic, it shows an improvement in power efficiency and data rate. Meanwhile, FDA is suitable for delay intolerant traffic, although it requires high power consumption compared to MDA since the location of the aggregator is fixed, and therefore the distance between the aggregator and the MTCD is not optimal. Therefore, it has been suggested to build a new approach that combines the two schemes into one scheme to satisfy the advantage of both; this third scheme is called CDA [18], where both fixed and mobile data aggregators cooperate to aggregate data from massive MTCDs (mMTCDs).

The FDA is assigned to aggregate data from delay intolerant mMTCDs, while the MDA is used to aggregate data of delay tolerable mMTCDs. The single point of failure and the suboptimal location of the FDA are avoided. A dynamic resource allocation based on the priority of MTCDs is presented. Although the results show good performance in term of outage probability, energy efficiency, and system capacity, resource allocation managed by the eNB and the aggregator play no role in resource assignment; it simply forwards the resource request from the MTCDs to the eNB. In particular, resource allocation is assigned based on the availability and the number of resources requested by MTCDs individually, which contradicts the concept of aggregation—that the resources blocks are assigned to the aggregator instead of individual assignment to each MTCD.

### 4) DATA AGGREGATION IN CAPILLARY NETWORK

Data aggregation in capillary networks connected to an LTE is introduced in [23], [26], [41]. In [23], fixed MTCDs are grouped to one fixed aggregator with a capillary connection; the aggregator is connected to the LTE BS by a cellular

channel. A fixed aggregation period is considered, which creates an increase in packet delay. The trade-offs between random access interaction, resource allocation, and communication latency are presented, and the results show a clear reduction in access interaction and resource allocation, at the cost of increasing the packet delay during transmission. Similar results are presented in [22], in which an experimental study is implemented to evaluate the impact of data aggregation on the signaling overhead and delay. The results show a significant reduction in the signaling with data aggregation; they also show that the signaling load reduction improves as the aggregation level increases (the number of aggregated MTCDs). The study also shows a trade-off between delay and aggregation level, since the aggregation level increases as the traffic delay increases. However, it does not provide any details about resource management or QoS differentiation for different types of MTCD traffic.

The work in [22] is expanded by [42]. The author proposes a priority-based data aggregation scheme for MTCD communication over the cellular network; three types of MTCD data traffic with different priorities based on their delay requirements are presented. The author also validates the study by introducing an analytical model for the aggregator using an M/G/1 queue. The study shows good performance in terms of average waiting time and system delay, but this comes at the cost of increasing the power consumption. However, this study does not address the issues related to LTE resource allocation or MTCD traffic modeling. In addition, the study supposes that the MTCD traffic has a higher priority than H2H in the case where they approach their tolerable delay threshold; therefore, in the event of a high MTCD traffic rate, the improvement in the performance of MTCD will be at the cost of degrading the performance of H2H traffic.

In [26], the authors propose a group-based radio resource allocation model, in which MTCDs are grouped based on identical transmission protocols (such as WiFi, wireless personal area network (WPAN), ZigBee) and QoS requirements (data rate and delay) to ensure QoS levels for MTCDs. The authors take into consideration the following assumptions: the uplink of SC-FDMA based LTE-A networks, WiFi grouping for MTCDs, and common service features of MTCDs. They utilize an effective capacity concept to model a wireless channel in terms of QoS metrics. The authors formulate a framework as a sum-throughput maximization problem, which satisfies all the constraints associated with SC-FDMA RBs and power allocation in LTE-A uplink networks. They solve the resource allocation problem by transforming it into a binary integer programming problem and then formulate a dual problem using Lagrange duality theory.

In [41], an energy harvesting gateway is proposed as an aggregator, which is connected to the eNB through an LTE interface, while it is connected to MTCD through a capillary communication technology such as ZigBee IEEE 802.15.4. SC-FDMA resource allocation is studied, and the performance of the system in terms of data transmitted, the number of RBs, and the drop rate is evaluated. The evaluation of

the system is expressed as an optimization non-deterministic polynomial-time (NP-hardness) problem, and two transforms are applied to express the problem in a linearly separable format. A heuristic algorithm for resource allocation is also introduced and compared to the optimization solution. The data energy causality, delay, and SC-FDMA constraints are taken into consideration. TABLE 1 summarizes the comparison between some data aggregation studies in the literature.

### C. DATA AGGREGATION BASED ON SDN

The rapid increase of data traffic in the core network, requires data aggregation to improve the performance of the system, in particular, for balancing the link loads. An aggregation approach with an admission control to provide a QoS for the SDN is introduced in [43]. The authors suggest rejecting the incoming data flows if it causes a degrading in the performance of the already admitted flows. Based on the performance metrics of the already admitted flows, the SDN controller is used to take a decision of accepting or rejecting the incoming flows. This study shows a reduction in packet loss ratio and delay.

An aggregation and scheduling approach in the flow-level for smart metering is proposed in [44]. The authors focus on investigating the fairness for traffic flows using SDN's flow-level features. Although the flow aggregation proposed for smart metering improves the overall throughput of the system, it experiences a problem of unfairness. So the authors used NS-3 and Mininet based evaluation to prove that their aggregation and scheduling approach achieves fairness for smart meters.

An efficient approach of flow aggregation for the delay-insensitive traffic control based on SDN framework is proposed in [45]. The study focuses on the case of massive number of small delay-insensitive traffic flows. The authors introduced a new data structure called flow tree, which is used to aggregate and decompress traffic flows according to the flow size in such a way to be adaptive to the changes in network conditions. This approach reduces the cost of communication between the controller and OpenFlow switches, and the cost of storage in switches memory.

Due to the expected increase in the data traffic from a huge number of sensors used in IoT applications, and given that the header in IoT packets consumes a large percentage of the total packet's size, this causes high overhead. The data aggregation based on SDN was one of the effective solutions to reduce the message delivered to the SDN controller. An aggregation/disaggregation approach based on SDN has been introduced for data sensors in IoT applications [46]. The authors exploited the (P4) switches proposed in [47]. Two P4 switches were used. One switch was used for receiving all data packets from IoT sensors, buffering them, and concatenating them with some metadata into a large packet transferred to the second P4 switch. The second switch in turn performs disaggregation to extract the original packets. A noticeable delay is shown in the process of disaggregation.

The authors analyzed their work using IoT talk platform. They showed a decrease in packet loss, improvement in system throughput, and reduction in communication between the SDN controller and switches.

The work presented in [46] introduces a mathematical analysis of the generated streams from the gathered packets, without including the designing and implementation issues, and without reporting the maximum throughputs. A similar work in [48] proposed by the same authors, involves implementation and design issues related to the aggregation and disaggregation approaches and their measured throughputs. The results show an improvement in the maximum throughput during aggregation, but a noticeable delay was incurred during disaggregation process. Moreover, they extended their work in [49] by solving the limitation of fixed payload size and the maximum number of aggregated packets, by supporting different payload sizes and allowing any number of aggregated packets as long as it does not exceed the maximum transmission unit (MTU). In addition, the aggregation and disaggregation throughputs were improved and can reach the line rate (i.e., 100 Gbps).

The authors in [50] proposed a second layer (L2) communication protocol for the Internet of Things programmable data planes referred to as Internet of Things Protocol (IoTP). The main goal of this protocol is to achieve the data aggregation algorithms within the hardware switches, at the network level. This process takes into consideration the network status and information such as MTU, delays, link bandwidths, and underlying communication technology, to enable the data aggregation algorithms dynamically. It provides support for different IoT communication technologies, different aggregation algorithms, and implementations of multi-level data aggregation. They implemented IoTP based on P4 language and using emulation-based Mininet environment. They showed a noticeable improvement in data aggregation.

In [51], the authors proposed an LTE-WiFi spectrum aggregation (LWA) based on the M-CORD platform which is used as an SDN platform to provide network function virtualization (NFV), cloud computing, edge computing, and virtualized RAN capabilities. They integrated WiFi with LTE in a very tight coupling scheme. Data from both networks is aggregated at the LTE PDCP layer, while a top-level network configuration is supported to the network orchestrator (XOS) of the M-CORD. They showed a significant improvement in system throughput compared to other similar scenarios. The traffic was split between LTE and WiFi based on the packet number: the even number are sent to LTE and the odd numbers to WiFi. This reordering function caused an increase in the packet delay.

In [52], the authors proposed an LTE-WiFi data aggregation on the RAN level based on the assistance of SDN (LWA-SA). They supposed a dual connectivity UE to both LTE and WiFi. Traffic was then split between LTE and WiFi based on the QoS requirements, and the best WiFi access point (AP) was elected using a Genetic algorithm (GA). SDN

**TABLE 1. State of the art comparison of different studies on data aggregation.**

Ref. No.	Type of data Aggregator (type / sub-type)	Mobility of MTCDs?	Include QoS?	Analytical tool	Performance metrics	Advantage	drawback
[23]	Capillary based / single FDA	No	No	- Simulation and Poisson random process	- Average delay - Normalized number of PRB required.	Capillary based data aggregation	- QoS has not been supported. - Fixed aggregation period leads to large packet delay.
[22]	Capillary based / single FDA	No	No	- Experimental study - M/G/1 queuing model	- Avg. packet delay - Signalling overhead	- Capillary based data aggregation - Experimental measurements used to proof the study	- QoS has not been supported. - Resources allocation has not considered.
[9]	Capillary based / single FDA	No	Yes	- M/G/1 queuing model	- power consumption - packet delay	- QoS provisioning. - A trade-off between power consumption and packet delay is provided.	- Adaptive level of aggregation has not included. - Max aggregation delay for each class dose included
[27]	Cellular based / single FDA.	Yes	No	- Simulation	- Mean PRBs utilization - End-to-End delay	shows PRB efficiency	- QoS has not been supported, - Fixed number of PRBs assigned to RN
[28]	single FDA	No	No	- Stochastic Geometry	- Energy consumption - Rate coverage - Outage rate - SIR Vs. Num. of hops (aggregation steps)	Optimal Number of hops Vs. energy consumption is investigated	- QoS does not be supported
[18]	Cellular based / Mobile multiple CDA	No	Yes	- Outage Probability	- Outage probability. - Energy efficiency. - system capacity	- Cooperative Data aggregation is proposed. - Delay tolerant, and delay intolerant traffics are included in the study	The utilization of RBs and multiplexing are not considered
[30]	Cellular based /multiple FDA	No	Yes	- Outage probability with Stochastic Geometry	- Outage probability. - Transmission capacity.	- Location based relaying. - Overcome packet duplicated. - the dynamic resource allocation scheme	- Fixed number of PRBs assigned to the aggregator. - Only consider a homogeneous type of MTCDs. QoS has not been supported
[29]	Cellular based / single FDA	No or Low mobility	No	- Outage probability with Stochastic Geometry	- Probability of successful channel utilization. - MTCD success probability	Data aggregation and resource allocation in large scale network have been investigated.	- QoS has not been supported, - Fixed number of PRBs assigned to RN
[17]	Cellular based / Single MDA (UE)	No	No	- TDMA slots reservation is represented as a "balls into bins problem" with a Markov Process.	- The average number of MTCDs served. - average transmit power per served device - delay	- Exploiting D2D connectivity to aggregate MTCD. -	- QoS not considered
[42]	Cellular based / Single MDA (UE)	NO	NO	- stochastic geometry with a Poisson Hard Sphere (PHS) Model	- SIR coverage probability. - probability of successful data delivery	- A large-scale MTCD data Aggregation by using D2D protocol.	- QoS not considered. - Only homogenous MTCD used.
[37]	Cellular based / multiple MDA (UE)	NO	NO	- optimization problem	- energy-efficiency - network lifetime	- Proposed an energy-efficient radio resource allocation (RRA) scheme for MTCD. - Data aggregation based on residual energy of MTCD	- QoS not considered.
[44]	Cellular based / multiple MDA (UE)	No	No	- a stochastic geometry	- Coverage probability. - Average data rate	- The mobility of UE is addressed. - Multi-hope relaying is introduced. - Improve data rate and coverage area.	Data aggregation has not implemented, but instead, they suggest using a UE to relay the data of MTCD in a multi-hope approach.
[63]	Cellular based / multiple MDA (UE)	Yes	No	- No analytical used. - No simulation used	NAN	- Three hops D2D based resource allocation for MTC is proposed.	- No Data aggregation is implemented, but multi-hope D2D relaying for MTC is proposed. - QoS not supported.
[32]	Cellular based / single MDA (UAVs)	No	No	- Simulation - Optimization method	- Power consumption. - Min. number of Aggregators. - Min. number of PRB per aggregator.	- Exploit drones as an RN for data aggregation.	- QoS not supported - Required extra protocols to compute the positions of MTCD and H2H . And to compute the Dwell time for UAVs over CH.
[40]	Cellular based / multiple MDA	yes	No	- Simulation - Probability of Random process	- Outage Probability - Transmission capacity	- Dynamic resource allocation and grouping for MTCDs. - Load balancing for each group.	- QoS has not been included. - Signalling overhead for location information exchange between MTCD and BS.

platform provides an intermediate layer for UE to aggregate the WiFi traffic through LTE without the need for an interface Xw between LTE and WLAN. They used a Lagrange Multiplier Method to compute the throughput maximization as an optimization problem that satisfies the constraints of power and interference.

A novel SDN based smart gateway (Sm-GW) was introduced in [53]. A Sm-GW was inserted between small cell eNBs and the multiple operators' gateways such as LTE S/P-GWs. In order to manage the backhaul link capacity, a scheduling algorithm was suggested for backhaul resource sharing with the assistance of SDN orchestrator. The results showed that SDN orchestrator provided flexible resource

management between the Sm-GWs, and hence improved the utilization of the backhaul bandwidth.

A Fog computing based Sm-GW for IoT e-Health application was presented in [54]. The proposed system exploits its position between the LAN/PAN/BAN and WAN to collect health and context information from different sensors. It included different services such as local data processing, local storage, data mining, data security and privacy, in addition, to data transmission controlling, enabling efficiency in term of energy and communication bandwidth. An intelligent intermediate layer was introduced between sensor nodes and the cloud to provide smooth and efficient e-Healthcare services while supporting patients' mobility. Complete system



implementation was presented, in addition to an Early Warning Score (EWS) notification system to inform for any emergency case.

In [55], a gateway for the Cloud of things (CoT) was introduced for managing things and to represent data for the end user. The lightweight virtualization technologies were exploited to improve the efficiency of the designed gateway and to decrease the impact on the performance. It mitigated the unnecessary communications between the gateway and Things, and therefore, reduced energy consumption. However, this study has some limitations, as it needs more adaptation algorithms to reduce the communication between things and the cloud [56].

Fog Computing platforms with Sm-GW has been proposed for IoT devices and wireless sensors in [57]. The main purpose of Fog Computing is to insert an intermediate layer between underlying devices and cloud network to provide preprocessing, monitoring, storage, and security. The Sm-GW plays an important role in achieving these functions. Furthermore, Sm-GW used to filter and mitigate the IoT communications by performing data pruning before sending them to the cloud server while meeting the constraints of the underlying devices and satisfying the requirement of high-level applications.

A Sm-GW based on Fog Computing was proposed in [58]. It has the ability to analyze the data before transmitting it to the cloud, and can differentiate between real-time data and non-real time data. Thus in order to utilize the available bandwidth efficiently, it responds to real-time data and sends it to the cloud directly, while the non-real time data is pre-processed, filtered and only the meaningful data is sent to the cloud.

### III. THE CONTRIBUTION OF THIS PAPER

This paper introduces a QoS-based data aggregation algorithm for MTCDs and resource allocation in an LTE-A network. Various types of MTCDs with different QoS requirements are considered. An aggregator is designed inside the RN (a layer-3 in-band LTE-A RN) to aggregate data from different types of MTCDs, process it, reformat it, and then relay it to the LTE eNB. The processing task consists of classifying the data into three priority classes, then buffering it so as not to exceed its delay tolerance threshold, and then sending it to the LTE eNB. The priority for each class is assigned based on its level of tolerance to delay. Unlike previous research, this paper uses an adaptive maximum aggregation delay and an adaptive Transport block size (TBS) threshold. These two parameters are very important for controlling the aggregation process to increase resource utilization efficiency with a minimal cost of delay.

Two resource allocation and scheduling schemes are used in this paper, a data buffer aware scheduling scheme and a moving boundary point scheme. In the former, the LTE resources are partitioned between aggregated users (MTCDs connected to RNs) and regular users (H2H) in a proportional approach to their data buffer size. In the second scheme,

the LTE resources are shared and partitioned in a hybrid manner to guarantee a minimum requirement RB for H2H, while also preventing the MTCDs from entering a starvation state.

A simulation model using MATLAB is designed to analyze the system performance in terms of throughput, utilization, loss ratio, and average packet delay. This paper also presents a survey of the literature covering the majority of works in the field of study, including smart gateway, data aggregation based on SDN and Fog Computing and the new works in 2020 (i.e., [41], [40], [50]–[53], [58]).

### IV. SYSTEM MODEL

The system considered to evaluate the proposed algorithm is as shown in FIGURE 1. It consists of one LTE Base Station eNB a number of MTCDs, coexisting with a number of H2H devices supported by LTE. The H2H devices are assumed to connect directly to LTE BS, while the MTCDs are first connected to the RN acting as aggregator then to the LTE BS. Three fixed Layer-3 In-band RNs characterized according to the 3GPP specifications in [59] are installed within the coverage area of LTE BS. Each RN works as an intermediate node to serve the MTCDs within its coverage area. Each RN has dual interfaces and dual functions: works as a base station from the point of view of users through Uu interface, and as UE from the point of view of LTE base station through Un interface.

In order to efficiently manage the resources of the LTE eNB base station and the RN, we assumed that the MTCDs within the coverage area of each RN are clustered and aggregated using an aggregator implemented inside the RN as shown in FIGURE 2. The aggregator collects the packets sent by the MTCDs connected to the RN, may delay them, reformats them and forwards them to eNB. Through this process, the small packets from MTCD are aggregated and reformatted such that the LTE RB assigned to RN can be exploited more efficiently since a single RB has more capacity than what is needed by one MTCD.

The aggregator is added since one RB allocated by eNB to MTCD provides more capacity than what may be needed by a single MTCD. Thus packets generated by MTCDs are first aggregated, as shown in FIGURE 3, and then allocated

TABLE 2. Characteristics of MTCD and H2H traffics considered in this work.

		Packets size	Distribution type	Mean inter-Arrival rate
MTCD Traffic type	e-Healthcare	32 B	Poisson	15 TTI
	Traffic monitoring	150 B	Poisson	20 TTI
	Smart metering	200 B	Poisson	30 TTI
H2H Video		100 – 200 B per Slice	Pareto	2.5 – 12 TTI
	TTI : Transmission Time Interval			

RBs according to some policies that will be defined later on. This technique will be very efficient for MTC applications generating small packets that are delay tolerant.

Three types of MTC traffic sources, namely: e-Healthcare, traffic monitoring and smart metering, and one H2H application with video traffic will be considered. The MTC sources will be served according to semi-priority scheme, where the e-Healthcare source will have the highest priority, while the smart metering will have the lowest. The traffic source characteristics of each one of the four types is as shown in TABLE 2.

**V. RESOURCE MANAGEMENT IN THE PROPOSED SCHEMES**

The resources in eNB are allocated to UEs and MTCs in two stages. In the first stage, the eNB PRBs are partitioned between direct UEs and RNs, the RBs assigned to RN are exploited to transmit the data buffered within aggregator to eNB through the backhaul link (link between RN and eNB), while in the second stage, the active MTCs reuse the sub-carrier that doesn't used in backhaul link to transmit their data to RN, in such to avoid the self-interference. We assumed there is no interference between the access links from one RN cluster to another RN cluster.

Two schemes were used for partitioning the resources between regular users UEs and RNs. In the first scheme, denoted here by the proportional fairness, the LTE resources are partitioned between regular users and RNs proportionally to the data buffered in each one. In the second, however, a moving boundary point is used to split LTE resources into three. One part is reserved for MTCs, a second for H2H users, and the third is shared between H2H and MTCs according to their requirements. A hard threshold value is used to partition the resources between H2H users and MTCs. The next sub-sections explain the two schemes of resource partitioning.

**A. PROPORTIONAL FAIRNESS RESOURCE PARTITIONING SCHEME**

In the proportional fairness scheme, the LTE resources are partitioned between H2H UE users, and MTC based on their buffered data size. The proposed resource allocation scheme is implemented in two stages. In the first stage, the LTE resources are partitioned between H2H direct users and the backhaul link of the relay users based on the size of data buffered at each H2H and each RN respectively. A buffer aware proportional fairness algorithm is used, it is similar to the algorithm in [60]. While resources are partitioned in [60] between the RN and direct users based on the number of users attached to each RN and the number of direct users, in our proposed algorithm the resources are partitioned based on the data buffered at the aggregator inside the RN and the data buffered at each regular user.

The Buffer State Report is used to inform the eNB about the amount of data buffered within its clients, then the LTE resources are assigned to RNs and UE according to the

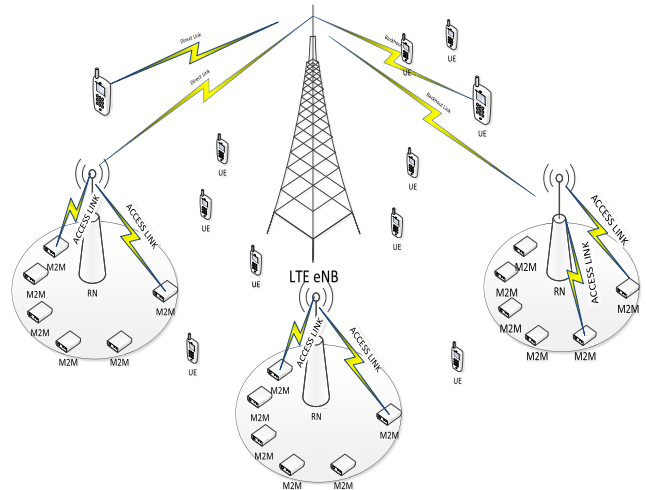


FIGURE 1. System model.

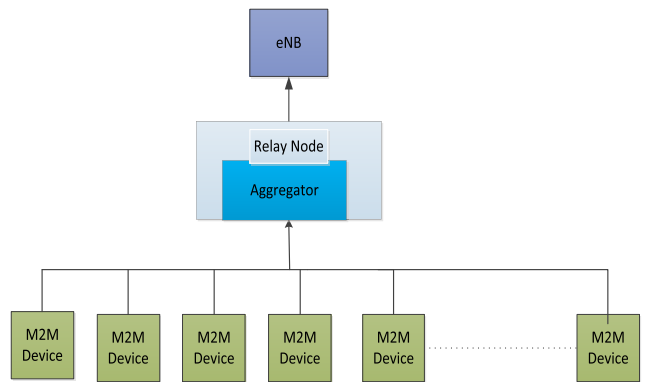


FIGURE 2. Hierarchical Architecture of the System model.

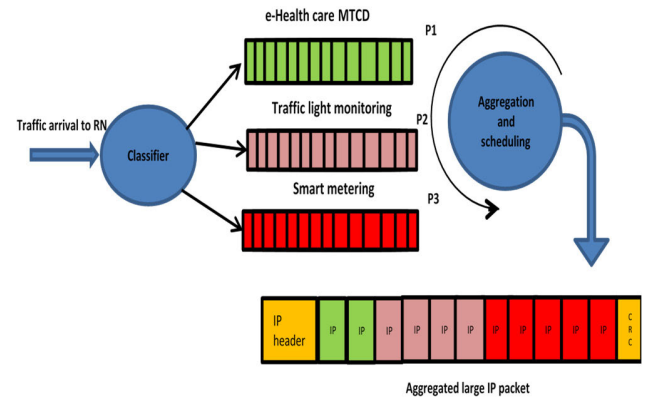


FIGURE 3. A QoS based Aggregator Scheme.

following equations:

$$RB_{RNj} = RB_{tot} \times \frac{BF_{RNj}}{\sum_{i=1}^N BF_{RNi} + \sum_{m=1}^H BF_{UEm}} \quad (1)$$

$$RB_{UE} = RB_{tot} \times \frac{\sum_{m=1}^H BF_{UEm}}{\sum_{i=1}^N BF_{RNi} + \sum_{m=1}^H BF_{UEm}} \quad (2)$$

where  $RB_{tot}$  is the total number of Resources Blocks in one LTE sub-frame;  $BF_{RNj}$  is the size of data buffered in the  $j$ th RN;  $BF_{UEm}$  is the size of data buffered in the  $m$ th H2H

user;  $N$  is the number of RNs;  $H$  is the total number of H2H devices;  $RB_{RNj}$  is the number of RBs assigned to the  $j$ th RN;  $RB_{UE}$  is the portion of RBs assigned to all regular users. The  $RB_{UE}$  is distributed to all regular users (H2H users) in a Round Robin manner. RN exploits the resources that have been granted by eNB to send the aggregated packets at the aggregator's queue through the backhaul link to eNB, by serving the aggregator's queues based on their priority, starting to serve the high priority queue which has buffered the e-Healthcare traffic, then serving road monitoring traffic, finally serve the smart metering traffic. An additional improvement is added to provide a balance between the second priority and third priority in the case where the delay of Head of Line (HOL) packet at the third queue reaches the threshold value, while the HOL packet at the second queue has a tolerable delay, more explanation in section VI.

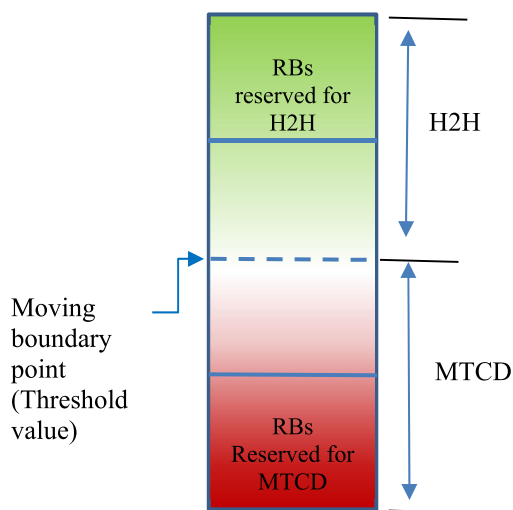
In the second stage, the RBs for MTCDs connected to RN in the access link are assigned, we suggest that MTCDs use LTE SC-FDMA, the RN manages the available resources by reusing the frequencies used by other RNs. In more details, we suggest that the RN are spatially isolated, where the RBs used in the access links of RN can be reused by the other RNs while avoiding self-interference between  $U_u$  and  $U_n$  interfaces of RN; the RBs used by RN in backhaul link cannot be used in access link in the same TTI. For simplicity, we suggest that the RN uses a round-robin mechanism to manage the available resource and allocate them to the active MTCDs in their transmission in the access links. Whenever the MTCD acquires a resource from RN, they use them to send their data to RN, and then the aggregation function takes place and implemented within RN.

**B. MOVING BOUNDARY POINT RESOURCE PARTITIONING SCHEME**

This scheme provides a hybrid mechanism for resource sharing and partitioning between H2H and MTCD users. First, the number of RBs required by each user is estimated based on the size of data buffered for each user, and the channel quality indicator (CQI) for each user. In the same approach, the number of RBs required by each RN is estimated based on the data buffered at each RN and the CQI between RN and eNB. The RBs in this scheme are divided into three parts. The first part is reserved for H2H to guarantee the minimum data rate for each H2H user, known as the guaranteed bit rate. The second part is reserved for MTCDs to guarantee RBs for high priority MTCD traffic, while the third part is shared between the H2H users and MTCDs. The shared part of RBs can be exploited by any type of user based on their requirement to ensure that the delay tolerance value is not exceeded. A predefined moving boundary point is set as a threshold value to split the shared part of RBs between H2H and MTCD users; this threshold value is elastic and can be varied to increase the RBs assigned to H2H users in the event that there are free RBs in the other part, and vice versa.

In the event that there is only one type of user requiring a resource block, while there is no data awaiting transmission

by the other type of user, all RBs are available for the type of users that need the RBs. FIGURE 4 shows the concept of moving boundary point resource sharing and partitioning. This scheme guarantees that some RBs are reserved for H2H, and therefore guarantees that H2H users are not affected by the huge number of MTCDs. At the same time, it avoids MTCDs entering a starvation state and guarantees at least meeting a minimum RBs allocation level for high priority MTCDs. In addition, it provides elastic resource partitioning between H2H and MTCD users. The moving boundary point can be adjusted to increase the resources assigned to H2H, but this comes at the cost of RB assignment to MTCDs.



**FIGURE 4. LTE RB partitioning and sharing based on moving boundary point scheme.**

**VI. AGGREGATION FUNCTION**

The aggregation function is implemented within the RN, and it takes place when MTCD traffic arrives at the RN. The function aggregates all data from different types of MTCD and classifies them into different queue buffers based on their priority. The size of each queue inside the aggregator is assumed to have infinite capacity, and when the delay of aggregated packets exceeds their tolerable delay limit, they will be dropped. Each class has its own buffer in the RN, as shown in FIGURE 3. We assume three types of traffic with different priorities: the first (highest) priority for eHealth traffic, the second priority for MTCD traffic monitoring, and the third (lowest) priority for MTCD smart metering traffic.

The aggregator accumulates the traffic in its buffers from different MTCDs and delays them until it reaches one of two parameters: the maximum tolerable delay threshold  $D_{i,max}$  or the maximum buffer size threshold  $Buf_{max}$  (whichever occurs first). These two parameters are very important for controlling the aggregation process and they should be selected in such a way as to improve RB utilization while maintaining traffic delay at a level below the tolerable delay threshold of any traffic. As the aggregation delay increases, so the resource utilization increases up to a limit beyond which any increase of delay aggregation will degrade the system's

performance in terms of increasing delay without any gain of resource utilization. Three different tolerable delay threshold values are used; one value for each type of traffic. The delay threshold for the highest priority traffic is the least, while the delay threshold for the lowest priority traffic will be the largest one. It is, therefore, expected that the aggregation delay for the high priority traffic will be less, while the lowest priority traffic will be delayed more. Each packet has its own timer, and when it reaches the tolerance threshold value it triggers the RN to request RBs to transmit the aggregated data. In the same approach, the size of the aggregated data should be accumulated until it reaches the threshold value, after which any increase in the size of data aggregated will degrade system performance.

The two parameters are computed adaptively to the type of aggregated data, their tolerable delay, and the estimated TBS of each RN.

$$Buf_{max} = TBS_{RN} - RNUnoverhead \quad (3)$$

TBS is defined as the number of bits that can be transmitted based on the RBs used, the modulation rate used, and the code rate. In an LTE RN, TBS depends on the number of RBs assigned to the RN. The CQI between RN and eNB indicates the modulation rate used to transmit the aggregated data. The aggregator traces the history of RB assignment to the RN (e.g., the last 10–15 assignments) and uses this to estimate the RBs that can be assigned to the RN in the next time slot (TTI). The aggregator then estimates the TBS for the next slot, and this TBS is used as the threshold value to which the aggregated data is accumulated.

$$TBS = nPRBs \times nDatasybol \times modulationrate \times coderate - CRC \quad (4)$$

where  $nPRBs$  = the total number of RBs assigned to the RN,  $nDatasybol$  = the number of data symbols within RBs in one subframe ( $12 \times 7 \times 2$ ),  $RE$  = the resource elements used for synchronization,  $modulationrate$  = the modulation order based on the SINR and channel quality between RN and eNB, and  $CRC$  is the cyclic redundancy check (equal to 24 bits in LTE).

When the RN is granted an RB to transmit its traffic, the aggregator collects the traffic from the different buffers starting with the highest priority queue, then the second priority, and finally the third priority, until the accumulated data fills the available granted RB. In this way, traffic with the highest priority is transmitted first. To keep the drop rate for each queue as low as possible while buffering packets in the RN, we suggest another priority for the traffic with same type in the same queue based on its tolerable delay; the packet with the lower tolerable delay being served first. To further improve system performance we suggest that the traffic of the third priority class can be served before the traffic of the second priority in the event that the HOL packet of the third priority class reaches its delay threshold while the delay of the HOL packet of the second priority has a tolerable

delay and can tolerate further delay without exceeding its delay threshold. This will provide a little balancing between the second- and third-class priorities and avoid the third priority class from entering a starvation state, thus decreasing both delay and drop rate. Of course, this comes at the cost of a small increase of delay for the second priority class.

## VII. PERFORMANCE METRICS

This section sets out the performance metrics used to evaluate the proposed algorithm in this study: system utilization, average packet delay, average drop rate, and average throughput.

### A. SYSTEM UTILIZATION

System utilization is one of the most important key performance indicators used to evaluate these types of systems. It is expected that the proposed aggregator will improve the performance of the system in terms of utilization by exploiting the RBs assigned to the RN efficiently. System utilization is defined as the average percentage of TBS used to transmit the aggregated data. In particular, utilization is defined as the effective throughput or spectral efficiency bits per second per hertz. In LTE, the TBS refers to the Physical Layer PHY payload to be transmitted over the radio interface, which consists of the MAC packet plus a 24-bit CRC overhead.

The average utilization for each class priority is computed by averaging the throughput of all users belonging to that class over the maximum throughput of the system; the throughput for each user is defined as the total number of bytes transmitted correctly over the simulation time.

The maximum throughput of the system is computed assuming the ideal case of the BS where the channel quality is optimum, the highest modulation rate and code rate are used, and by using the total number of RBs available in the system. The maximum throughput of the system can be defined as the maximum number of bytes that can be transmitted over time in the ideal environment of that system.

$$Utilization_i = \frac{\sum_{u=1}^{M_i} Throughput(u)}{MaxThroughput} \quad (5)$$

$M_i$  = the number of MTCDs belonging to priority class  $i$ .

### B. PACKET DELAY

Packet delay is calculated from the time the packet is generated by the user until it arrives at the eNB. It comprises two terms of delay: the delay of the packet within the buffer of the user, and the delay of the packet inside the buffer of the aggregator. Packet delay is computed by creating a timestamp for each packet when it is generated. When the packet arrives at the eNB, the delay is calculated for each packet; after that, the average delay for all packets belonging to the same user is calculated. In addition, the average delay for all users belonging to the same priority class is computed, to compute the average delay for each priority class. Packet delay is an important metric, and it must not exceed the tolerable delay of each type of traffic. It can be used to evaluate how the proposed aggregator fulfills the QoS of each traffic type.

### C. PACKET LOSS RATIO

Due to holding packets within the user buffer until it is granted a chance to transmit its data directly to the eNB or to the aggregator, and the delay of the packets within the aggregator, some packets may exceed their tolerable delay, be dropped, and be considered as lost. We assume a dropped packet only occurs as a result of a delay exceeding the tolerable threshold and not due to an error in transmission. The loss ratio of each user is computed by dividing the number of lost packets by the number of packets generated by that user; the loss ratio is averaged for all users belonging to the same priority class to determine an average loss ratio for each class.

### D. AVERAGE THROUGHPUT

Average throughput is one of the major performance key indicators for evaluating communication systems. Throughput is defined as the amount of data transmitted by each user correctly over a given time period. In this paper, the amount of data transmitted by each user is computed in each TTI, and then, averaged over all users with the same priority class. This is then averaged over all simulation time slots to determine the average throughput for the simulation time. Throughput is used to measure and evaluate system provisioning of the QoS.

## VIII. SIMULATION MODELS

To evaluate the proposed algorithm, we have built a Matlab simulation model based on the one in [61], [62]. The model in [61] was modified to support LTE-A uplink transmission, and the RN was upgraded with a built-in aggregator. The simulation program was run and repeated twice with the same parameters shown in Table 3, once to evaluate the first proposed scheme (i.e., Proportional fairness resource partitioning scheme), while the second to evaluate the second

TABLE 3. Simulation parameters.

Parameter	Value
Cell radius	500 m
TTI time slot	1 ms
Simulation time in TTI	20000 TTI
System bandwidth	5 MHz
No. of BSs	1
No. of MTCD users	105 equally divided by type: 35 e-Healthcare, 35 road monitoring, 35 smart metering.
Packet size	e-Healthcare = 32 B (Bytes) Road monitoring devices = 150 B Smart metering = 200 B
No. of H2H users (LTE-UE)	45
No. of RNs	3
Type of RN	In-band layer 3 and fixed RN
No. of RBs assigned to RN	Adaptive to the size of their data buffer
Time threshold of the aggregator	Adaptive to the tolerable delay of each traffic = 5, 6, 50 TTI for e-Healthcare, road monitoring and smart metering traffic, respectively.
Transmission direction	Uplink
RN diameter (RN cluster range)	10 m
nRX, nTX for users	1, 1
nRX, nTX for RN	1, 1
Transmission mode	SISO
Noise figure	5 dB

proposed scheme (i.e., Moving boundary point resource partitioning Scheme).

In the first scheme, one LTE-A cell, and three types of MTCD with various traffic characteristics is supposed. The simulation was run by varying the mean inter-arrival rate for the MTCD, while fixing it for H2H users.

The mean arrival rate of MTCD traffic was increased in each run by increments of 5% of the initial value, where the initial values of mean arrival rate was  $\frac{1}{15}$  for e-Healthcare traffic,  $\frac{1}{20}$  for road monitoring traffic, and  $\frac{1}{30}$  for smart metering traffic. The smart metering traffic had the lowest arrival rate, while the road monitoring traffic had the second-highest arrival rate, and the e-Healthcare traffic had the highest arrival rate.

The simulation was run using the aggregator as described in Section VI, and the resource management as described in Section V. The simulation was also run without using the aggregator, allowing all MTCD and H2H users to connect directly to the eNB and using a round-robin scheduling algorithm, and then, the results were compared.

In the second scheme, the simulation was run with the same configuration parameters as in the first scheme, while also using the moving boundary point for resource sharing and partitioning between H2H and MTCD users as described in Section V-B. And the results were also compared to the case where the aggregation does not been used.

## IX. RESULTS AND PERFORMANCE ANALYSIS

This section is divided into two subsections, the first section presents the results for the first proposed schemes (i.e., proportional fairness scheme, while the second section presents the results for the second proposed scheme (i.e., Moving Boundary Point scheme).

All results in these two subsections will be presented as a function of the mean arrival rate. In each case, the simulation is run at 15 different mean arrival rates. The mean arrival rate is increased at each point with a 5% increment of the initial value of the arrival rate. Moreover, in each case, the simulation is repeated 10 times, and the simulation is run for 20000 TTI (i.e., 20,000 msec). The simulation result is averaged over the 10 runs with a 95% confidence interval. The observed simulation results a maximum error between runs of less than 0.50% of the mean value. The dashed curves represent the system performance without the aggregator, while the solid lines represent the system performance when using the aggregator.

### A. FIRST SCHEME: PROPORTIONAL FAIRNESS RESOURCE PARTITIONING (PFRP)

#### 1) UTILIZATION

The average utilization for all traffic types in both cases (with/without using the aggregator) is shown in FIGURE 5. In the case of using the aggregator (solid lines), the average utilization for all MTCD traffic increases as the mean arrival rate increases. This comes at the cost of a decrease in the

average utilization of H2H users (solid black line with “+” mark). The road monitoring traffic has the highest utilization, despite being medium priority; this is because the road monitoring traffic has the highest arrival rate.

By comparing the utilization in both cases (with/without using the aggregator), there is a significant improvement in the utilization for all MTCD traffic when using the aggregator: approximately 16% increase for the road monitoring traffic and approximately 2% increase for the e-Healthcare MTC traffic. Furthermore, the utilization improvement for video H2H traffic is close to 5%, although it decreases as the MTC arrival rate factor increases, because the resources are consumed by MTCD. In order to maintain H2H performance when MTCD traffic varies, a solution will be proposed in the second scenario.

A noticeable decrease in the utilization for smart metering traffic occurs at the high arrival rates. This is because the smart metering traffic has the lowest priority; therefore, the other types MTCD traffic is prioritized at the cost of performance deterioration of smart metering traffic.

Besides these differences in utilization among the various traffic types, the total utilization of the system with aggregator improves as traffic is increased, as shown in FIGURE 6.

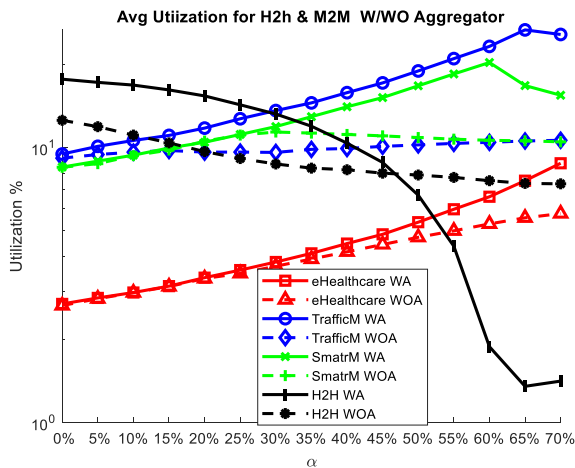


FIGURE 5. Average utilization for all type of traffic.

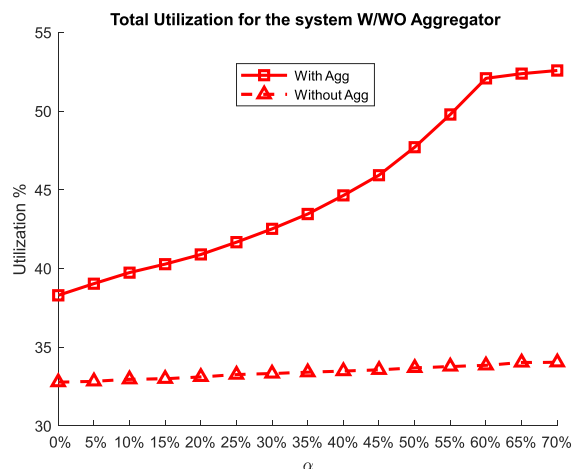


FIGURE 6. Total system utilization with/without using aggregator.

## 2) MTCD DELAY

In FIGURE 7, the average delay (in msec) for all types of MTCD traffic is presented on a logarithmic scale. In the case of using the aggregator, and at the low arrival rate factor, the results show a significant improvement in decreasing the average delay for all types of MTCD traffic, while at the high arrival rate factor, the average delay for smart metering traffic (solid green curve with “+” mark) and road monitoring traffic (solid blue curve with “○” mark) have a Higher average delay than in the case of not using the aggregator. This performance was expected, since the aggregator delays the traffic until a predefined data size or time aggregation level is reached. This comes against an increase of system utilization. The figure also shows that the smart metering traffic has the highest delay (solid green curve with “+” mark). This is because they have the lowest priority; therefore, they are delayed in the aggregator to provide higher performance for the highest priority traffic. The e-Healthcare traffic has the lowest average delay (since they have the highest priority); this validates the QoS provision of the proposed algorithm for all traffic in terms of delay.

## 3) PACKET LOSS RATIO

In FIGURE 8, the loss ratio of all traffic types is displayed on a logarithmic scale. In the case of using the aggregator, there is no loss ratio for the e-Healthcare traffic because they have the highest priority, while there is a slight loss ratio for road monitoring traffic (solid blue curve with “○” mark) at the high arrival rate factor. In addition, the figure shows a high loss ratio for smart metering traffic (solid green curve with “x” mark) at a high arrival rate. This is because the smart metering traffic has the lowest priority; therefore, at the high arrival rate, the system cannot serve all traffic and it starts to drop the traffic with the lowest priority. In case of not using the aggregator (dashed curves), the figure shows there is packet loss for all types of MTCD traffic; the loss ratio increases as the arrival rate factor increases, and their loss ratio is higher than in case of using an aggregator. The black color curves in the figure show the loss ratio for video

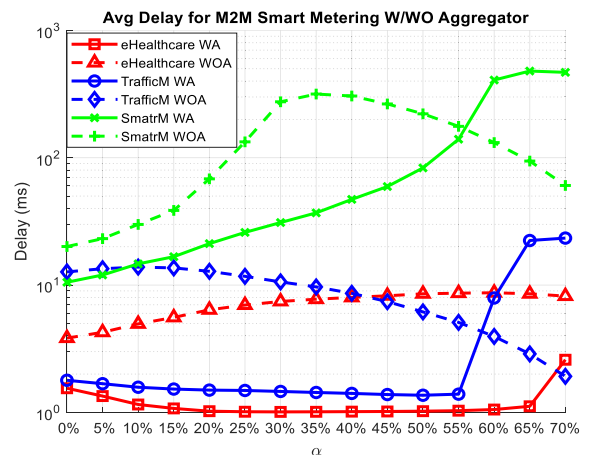


FIGURE 7. The average delay for MTCD traffic.

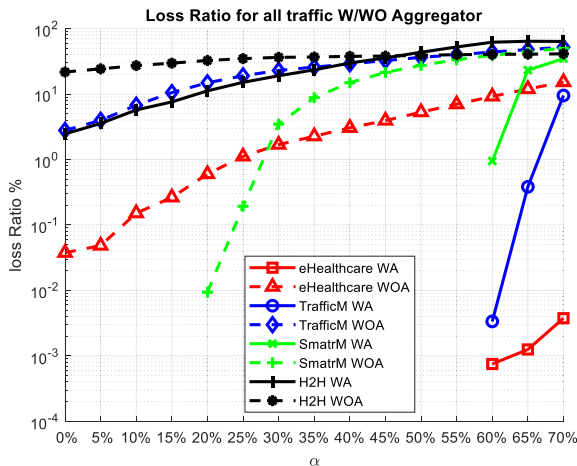


FIGURE 8. The loss ratio for all types of traffic.

H2H traffic; the loss ratio for video traffic in the case of using the aggregator (solid curve with “•” mark) is less than the loss ratio when not using the aggregator (solid curve with “|” mark) at the low traffic rate, while this result is reflected at the high traffic rate. This is because, at the high arrival rate factor of MTCD traffic, the MTCD consumes the resources; therefore, the resources allocated for H2H users decrease due to the resources being partitioned between MTCD and H2H users based on the data buffered on each one. This ultimately, leads to an increase in the loss ratio for H2H users.

In FIGURE 9, the loss ratio of the whole system is presented. It shows a significant improvement for the system in terms of decreasing the loss ratio by approximately 15% when using the aggregator at high traffic rates, while this improvement decreases to 6% at the low arrival rate.

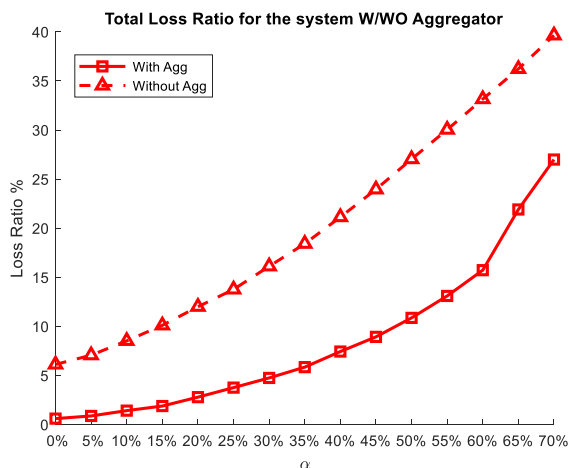


FIGURE 9. Loss Ratio for whole system.

4) TOTAL DATA TRANSMITTED

FIGURE 10 shows the total amount of data (in Megabytes MB) transmitted by MTCDs during the simulation time. It shows that the size of data transmitted increases as the arrival rate factor increases. The figure also shows that all MTCDs transmitted a larger number of bytes when using the

aggregator (solid curves), comparing to the data transmitted by the same devices without using the aggregator. At the high arrival rate factor, the smart metering traffic (green curve with “x” mark) shows a small decrease in data transmitted because they have the lowest priority, and there are insufficient RBs available for them. The road monitoring traffic transmits the largest volume of data, while the e-Healthcare transmits the least amount because of their data rate and packet size.

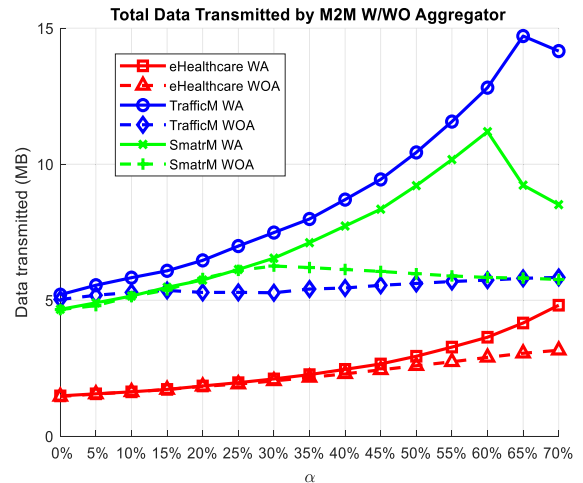


FIGURE 10. Total Data transmitted by MTCD device.

5) UTILIZATION IN TERM OF NUMBER OF RBs

FIGURE 11 shows the utilization of the system in terms of the number of resource block RBs used by MTCD and H2H users in both cases (with/without using an aggregator). It shows that using the aggregator decreases the number of RBs used by MTCDs (solid red line with “□” mark) against increasing them for the H2H regular users (blue curve with “○” mark). However, when the arrival rate of MTCD traffic exceeds a determined limit, the RBs used by the MTCD in case of using aggregator becomes greater than that used by the MTCDs without using the aggregator (at arrival rate 50%). By comparing the results of FIGURE 10 to the results of FIGURE 11, it is clear that MTCDs transmit a larger amount of data when using the aggregator, while using less RBs (solid red line with “□” mark).

B. SECOND SCHEME: MOVING BOUNDARY POINT RESOURCE PARTITIONING (MBPRP)

Although the PFRP scheme shows an improvement in the system performance in general, and for MTCDs in particular, it did not keep the performance of H2H users from the negative effects of increasing the MTC traffic. As was shown in the previous sections, the improvement of the MTCD comes at the cost of degrading the performance of H2H users. So a new scheme is proposed to provide a QoS for MTCD while maintaining the good performance of H2H users. The MBPRP scheme for resource partitioning between H2H users and M2M devices was described in section V-B, and its results are presented in next subsections.

1) UTILIZATION

Figure 12 presents the average utilization of all types of traffic in the MBPRP scheme, it shows that the utilization for H2H video traffic, e-Healthcare MTCD traffic, and road monitoring MTCD traffic and improved in case of using the aggregator. It also shows that this scheme keeps the utilization for H2H users at an approximately fixed level without or with a little effect of the increasing of MTCD traffic, but at the cost of the utilization of the lowest priority MTCD traffic (i.e., smart metering – solid green curve with ‘x’ mark) which is degraded as the arrival rate factor increase. Like the first scheme, the utilization of MTCD traffic with the highest, and the second-highest priority are increase as the arrival rate factor increases. This scheme guarantees that the highest MTCD traffic (e-Healthcare – solid red curve with ‘□’ mark) gets the required resources, at the same time keep the H2H users from the effect of arrival rate increasing of MTCDs.

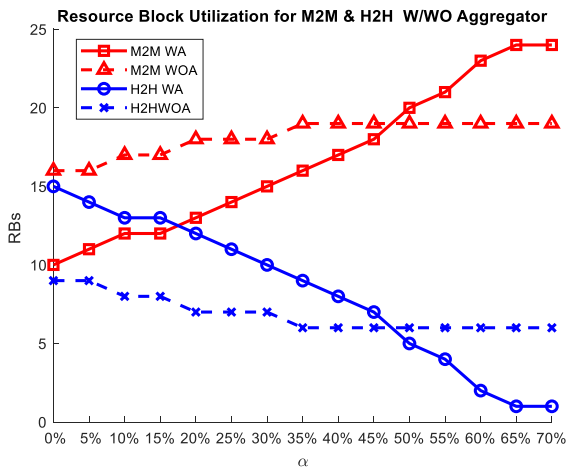


FIGURE 11. The Avg. Utilization of RBs for MTCD and H2H users.

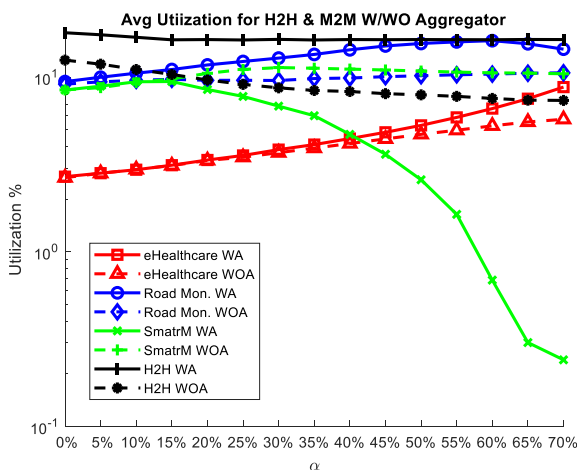


FIGURE 12. Avg. utilization for all type of traffic W/O Aggregator.

2) LOSS RATIO IN THE SECOND SCHEME

FIGURE 13 presents the total loss ratio for all types of traffic in the second scheme, it shows that the loss ratio increases as the arrival rate increases, it also shows that using the

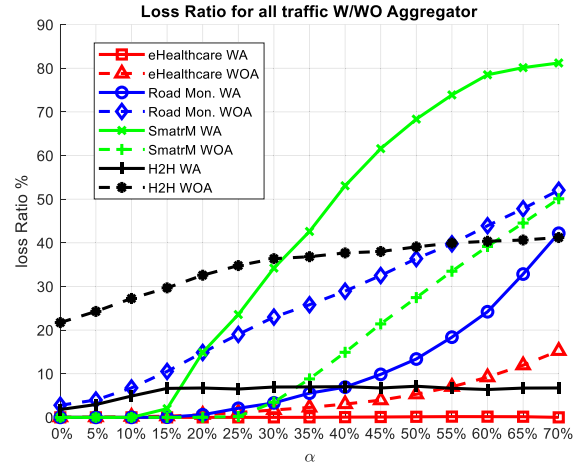


FIGURE 13. Total Loss Ratio for all types of traffic W/O aggregator.

aggregator (solid lines) decreases the loss ratio for all types of traffic except for smart metering traffic (solid green lines with ‘x’), this is because the smart metering traffic has the lowest priority. By comparing the result in FIGURE 13 to the results in FIGURE 8 it is clear that the loss ratio for H2H in the case of using aggregator is decreased in this scheme, the loss ratio for H2H traffic in this scheme does not exceed 7%. While the loss ratio for H2H traffic in the first scheme exceeds 60% as shown in FIGURE 8. This explains how this scheme keeps the H2H traffic from the negative effects of increasing M2M traffics. This comes at the cost of increasing the loss ratio of smart metering traffic.

X. CONCLUSION

A QoS based data aggregation algorithm was presented for the MTCDs traffic when integrated with the co-existent H2H users within LTE-A. The algorithm goal was to mitigate the effects of MTC traffic on the performance of H2H users while maintaining the QoS for each type of traffic. To achieve this, an aggregator with an adaptive aggregation delay for each type of traffic, and adaptive size of aggregation data has been used. Three types of MTCD traffic served with different priorities have been considered: e-Healthcare, road monitoring, and smart metering traffic.

Two resource allocation schemes have been presented: a proportional fairness data buffer aware resource partitioning and moving boundary point were considered. In the first, The LTE resources were partitioning between the RNs and H2H users proportionally to the size of data buffered, while In the second scheme, the LTE resources were partitioned and shared in a hybrid manner, by reserving some RBs for H2H to provide them with a Guaranteed Bit Rate GBR, and at the same time guaranteeing that the high priority M2M traffic does not get into starvation state.

The results showed a significant improvement in the system performance in terms of average utilization, number of resources used, loss ratio, and the average delay in the case where an aggregator was used. However, the first scheme has a limitation in isolating the H2H performance from



TABLE 4. Abbreviation.

Abbreviation	Definition
4G	4 <sup>th</sup> Generation
5G	5 <sup>th</sup> Generation
AP	Access Point
BAN	Body Area Networks
BS	Base Station
CDA	Cooperative data aggregation
CH	Cluster Head
CoT	Cloud of things
CQI	channel quality indicator
CRC	Cyclic Redundancy Check
CRS	Channel-aware resource scheduling
C-UEs	Cellular User Equipments
D2D	Device to Device
eNB	Evolved Node B
EWS	Early Warning Score
FDA	Fixed data aggregator
GA	Genetic algorithm
Gbps	Gigabits per second
GBR	Guaranteed Bit Rate
H2H	Human to Human
HOL	Head of the line
IoT	Internet of Things
IoTP	Internet of Things Protocol
LAN	Local Area networks
LTE	Long term evaluation
LTE-A	Long term evaluation-Advanced
LWA	LTE-WiFi aggregation
LWA-SA	LTE-WiFi Aggregation-SDN Assisted
M2M	Machine to Machine
MAC	Medium Access Control
MATLAB	Matrix Laboratory
MB	Megabytes
MBPRP	Moving Boundary Point Resource Partitioning
M-CORD	Mobile-Central Office Re-Architected as Datacenter (open source network lab)
MDA	Mobile data aggregator
MFDA	Multiple fixed data aggregators
MFDA <sub>s</sub>	Multiple fixed data aggregators
MHz	Mega Hertz
MMDA	Multiple Mobile Data Aggregator
mMTCDs	Massive machine type communication devices
MTC	Machine Type Communication
MTCD(s)	Machine Type Communication Device(s)
MTCGs	Machine Type Communication Gateway(s)
MTU	Maximum transmission unit
NFV	Network function virtualization
NP-hard	non-deterministic polynomial-time hardness
NRT	non-real-time
nRx	Number of Receivers
nTx	Number of Transmitters
OFDMA	Orthogonal frequency division multiple access
P4	Programming Protocol-independent Packet Processors
PAN	Personal area Networks
PDCP	Packet Data Convergence Protocol
PFRP	Proportional Fairness Resource Partitioning
PHY	Physical layer
PRB	Physical Resource Blocks
QoS	Quality of services
RACH	Random Access Channel
RBs	Resource Blocks
RE	Resource Elements
RN(s)	Relay Node(s)
RRM	Radio Resource Management

TABLE 4. (Continued.) Abbreviation.

RRS	random resource scheduling
RTTS	Remaining times to serve
SC-FDMA	Single Carrier Frequency Division Multiple Access
SDN	Software Defined Network
SFDA	Single fixed data aggregator
SINR	Signal-to-interference plus Noise Ratio
SIR	Signal to interference ratio
SISO	Single Input Single Output
SMDA	Single mobile data aggregator
Sm-GW	Smart gateway
TBS	Transport block size
TDMA	Time Division Multiple Access
TTI	Transmission Time Interval
UAV	unmanned aerial vehicle
UE	User Equipment
UL	Uplink
WiFi	Wireless Fidelity
WLAN	Wireless Local area Networks

the negative effects of increasing MTC traffic. This limitation was alleviated in the second scheme, where the results showed that the QoS for the H2H users was maintained while data rate of the MTCs was increased.

Although the proposed schemes provided significant improvements in system performance, the new trends in designing the data aggregator should exploit the new technologies such as SDN, fog computing and network virtualization, to design a smart gateway aggregator where the data analysis and resource allocation can be achieved with more flexibility. We suggest the researchers to combine our results with these new technologies to design more trusted and adaptive schemes in data aggregation.

## APPENDIX A

See Table 4.

## ACKNOWLEDGMENT

The authors thank the Deanship of Scientific Research and RSSU at King Saud University for their technical support.

## REFERENCES

- [1] S. Persia and L. Rea, "Next generation M2M cellular networks: LTE-MTC and NB-IoT capacity analysis for smart grids applications," in *Proc. AEIT Int. Annu. Conf. (AEIT)*, Capri, Italy, 2016, pp. 1–6, doi: 10.23919/AEIT.2016.7892789.
- [2] A. Whitmore, A. Agarwal, and L. Da Xu, "The Internet of Things—A survey of topics and trends," *Inf. Syst. Frontiers*, vol. 17, no. 2, pp. 261–274, Apr. 2015, doi: 10.1007/s10796-014-9489-2.
- [3] *Standardization of Machine-Type Communications*, document TR 0.2.4, 3GPP, Jun. 2014, pp. 1–12.
- [4] *System Improvements for Machine Type Communications (MTC)*, document TS, TR23.888, 3GPP, Release 11, pp. 1–34, 2012. [Online]. Available: <http://www.3gpp.org>
- [5] M. Z. Shafiq, L. Ji, A. X. Liu, J. Pang, and J. Wang, "A first look at cellular machine-to-machine traffic—Large scale measurement and characterization," *ACM SIGMETRICS Perform. Eval. Rev.*, vol. 40, no. 1, pp. 65–76, 2012, doi: 10.1145/2318857.2254767.
- [6] Cisco. (2020). *Cisco Annual Internet Report (2018–2023)*. [Online]. Available: <https://www.cisco.com/c/en/us/solutions/collateral/executive-perspectives/annual-internet-report/white-paper-c11-741490.html>
- [7] K. Zheng, F. Hu, W. Wang, W. Xiang, and M. Dohler, "Radio resource allocation in LTE-advanced cellular networks with M2M communications," *IEEE Commun. Mag.*, vol. 50, no. 7, pp. 184–192, Jul. 2012, doi: 10.1109/MCOM.2012.6231296.
- [8] A. Aijaz and A. H. Aghvami, "On radio resource allocation in LTE networks with machine-to-machine communications," in *Proc. IEEE 77th Veh. Technol. Conf.*, Jun. 2013, pp. 1–5, doi: 10.1109/VTC-Spring.2013.6692664.
- [9] P. Makris, D. N. Skoutas, N. Nomikos, D. Vouyioukas, and C. Skianis, "A context-aware backhaul management solution for combined H2H and M2M traffic," in *Proc. Int. Conf. Comput., Inf. Telecommun. Syst. (CITS)*, Athens, Greece, May 2013, pp. 1–5, doi: 10.1109/CITS.2013.6705719.
- [10] M. K. Giluka, N. Rajoria, A. C. Kulkarni, V. Sathya, and B. R. Tamma, "Class based dynamic priority scheduling for uplink to support M2M communications in LTE," in *Proc. IEEE World Forum Internet Things (WF-IoT)*, Mar. 2014, pp. 313–317, doi: 10.1109/WF-IoT.2014.6803179.
- [11] D. M. Soleymani, A. Puschmann, E. Roth-Mandutz, J. Mueckenheim, and A. Mitschele-Thiel, "A hierarchical radio resource management scheme for next generation cellular networks," in *Proc. IEEE Wireless Commun. Netw. Conf.*, Doha, Qatar, Apr. 2016, pp. 1–5, doi: 10.1109/WCNC.2016.7564706.
- [12] M. R. Mardani, S. Mohebi, and H. Bobarshad, "Robust uplink resource allocation in LTE networks with M2M devices as an infrastructure of Internet of Things," in *Proc. IEEE 4th Int. Conf. Future Internet Things Cloud (FiCloud)*, Vienna, Austria, Aug. 2016, pp. 186–193, doi: 10.1109/FiCloud.2016.34.
- [13] C. Ide, B. Dusza, and C. Wietfeld, "Performance of channel-aware M2M communications based on LTE network measurements," in *Proc. IEEE 24th Annu. Int. Symp. Pers., Indoor, Mobile Radio Commun. (PIMRC)*, Sep. 2013, pp. 1614–1618, doi: 10.1109/PIMRC.2013.6666400.
- [14] Y. Zhang, "Tree-based resource allocation for periodic cellular M2M communications," *IEEE Wireless Commun. Lett.*, vol. 3, no. 6, pp. 621–624, Dec. 2014, doi: 10.1109/LWC.2014.2366769.
- [15] Q. Li, Y. Ge, Y. Yang, Y. Zhu, W. Sun, and J. Li, "An energy efficient uplink scheduling and resource allocation for M2M communications in SC-FDMA based LTE-A networks," *Mobile Netw. Appl.*, pp. 124–140, 2019, doi: 10.1007/s11036-019-01400-w.
- [16] A. Aijaz, M. Tshangini, M. R. Nakhai, X. Chu, and A.-H. Aghvami, "Energy-efficient uplink resource allocation in LTE networks with M2M/H2H co-existence under statistical QoS guarantees," *IEEE Trans. Commun.*, vol. 62, no. 7, pp. 2353–2365, Jul. 2014, doi: 10.1109/TCOMM.2014.2328338.

- [17] G. Rigazzi, N. K. Pratas, P. Popovski, and R. Fantacci, "Aggregation and trunking of M2M traffic via D2D connections," in *Proc. IEEE Int. Conf. Commun. (ICC)*, London, U.K., Jun. 2015, pp. 2973–2978, doi: [10.1109/ICC.2015.7248779](https://doi.org/10.1109/ICC.2015.7248779).
- [18] T. Salam, W. U. Rehman, and X. Tao, "Cooperative data aggregation and dynamic resource allocation for massive machine type communication," *IEEE Access*, vol. 6, pp. 4145–4158, 2018, doi: [10.1109/ACCESS.2018.2791577](https://doi.org/10.1109/ACCESS.2018.2791577).
- [19] Z. Dawy, W. Saad, A. Ghosh, J. G. Andrews, and E. Yaacoub, "Toward massive machine type cellular communications," *IEEE Wireless Commun.*, vol. 24, no. 1, pp. 120–128, Feb. 2017, doi: [10.1109/MWC.2016.1500284WC](https://doi.org/10.1109/MWC.2016.1500284WC).
- [20] A. Lo, Y. Law, and M. Jacobsson, "A cellular-centric service architecture for machine-to-machine (M2M) communications," *IEEE Wireless Commun.*, vol. 20, no. 5, pp. 143–151, Oct. 2013.
- [21] N. Kouzayha, M. Jaber, and Z. Dawy, "Measurement-based signaling management strategies for cellular IoT," *IEEE Internet Things J.*, vol. 4, no. 5, pp. 1434–1444, Oct. 2017, doi: [10.1109/JIOT.2017.2736528](https://doi.org/10.1109/JIOT.2017.2736528).
- [22] N. Kouzayha, M. Jaber, and Z. Dawy, "M2M data aggregation over cellular networks: Signaling-delay trade-offs," in *Proc. IEEE Globecom Workshops (GC Wkshps)*, Austin, TX, USA, Dec. 2016, pp. 1155–1160, doi: [10.1109/glocomw.2014.7570073](https://doi.org/10.1109/glocomw.2014.7570073).
- [23] H. Shariatmadari, P. Osti, S. Iraji, and R. Jantti, "Data aggregation in capillary networks for machine-to-machine communications," in *Proc. IEEE 26th Annu. Int. Symp. Pers., Indoor, Mobile Radio Commun. (PIMRC)*, Aug. 2015, pp. 2277–2282, doi: [10.1109/PIMRC.2015.7343677](https://doi.org/10.1109/PIMRC.2015.7343677).
- [24] A. Laya, L. Alonso, and J. Alonso-Zarate, "Is the random access channel of LTE and LTE-A suitable for M2M communications? A survey of alternatives," *IEEE Commun. Surveys Tuts.*, vol. 16, no. 1, pp. 4–16, Feb. 2014, doi: [10.1109/SURV.2013.111313.00244](https://doi.org/10.1109/SURV.2013.111313.00244).
- [25] A. Aijaz and A. H. Aghvami, "Cognitive machine-to-machine communications for Internet-of-Things: A protocol stack perspective," *IEEE Internet Things J.*, vol. 2, no. 2, pp. 103–112, Apr. 2015, doi: [10.1109/JIOT.2015.2390775](https://doi.org/10.1109/JIOT.2015.2390775).
- [26] F. Ghavimi, Y.-W. Lu, and H.-H. Chen, "Uplink scheduling and power allocation for M2M communications in SC-FDMA-based LTE-A networks with QoS guarantees," *IEEE Trans. Veh. Technol.*, vol. 66, no. 7, pp. 6160–6170, Jul. 2017, doi: [10.1109/TVT.2016.2635262](https://doi.org/10.1109/TVT.2016.2635262).
- [27] S. N. K. Marwat, Y. Mehmood, C. Görg, and A. Timm-Giel, "Data aggregation of mobile M2M traffic in relay enhanced LTE-A networks," *EURASIP J. Wireless Commun. Netw.*, vol. 2016, no. 1, p. 14, Dec. 2016, doi: [10.1186/s13638-016-0598-0](https://doi.org/10.1186/s13638-016-0598-0).
- [28] D. Malak, H. S. Dhillon, and J. G. Andrews, "Optimizing data aggregation for uplink machine-to-machine communication networks," *IEEE Trans. Commun.*, vol. 64, no. 3, pp. 1274–1290, Mar. 2016, doi: [10.1109/TCOMM.2016.2517073](https://doi.org/10.1109/TCOMM.2016.2517073).
- [29] J. Guo, S. Durrani, X. Zhou, and H. Yanikomeroğlu, "Massive machine type communication with data aggregation and resource scheduling," *IEEE Trans. Commun.*, vol. 65, no. 9, pp. 4012–4026, Sep. 2017, doi: [10.1109/TCOMM.2017.2710185](https://doi.org/10.1109/TCOMM.2017.2710185).
- [30] U. Tefek and T. J. Lim, "Relaying and radio resource partitioning for machine-type communications in cellular networks," *IEEE Trans. Wireless Commun.*, vol. 16, no. 2, pp. 1344–1356, Feb. 2017, doi: [10.1109/TWC.2016.2645688](https://doi.org/10.1109/TWC.2016.2645688).
- [31] M. N. Soorki, M. Mozaffari, W. Saad, M. H. Manshaei, and H. Saidi, "Resource allocation for machine-to-machine communications with unmanned aerial vehicles," in *Proc. IEEE Globecom Workshops (GC Wkshps)*, Dec. 2016, pp. 1–6, doi: [10.1109/GLOCOMW.2016.7849026](https://doi.org/10.1109/GLOCOMW.2016.7849026).
- [32] M. Mozaffari, W. Saad, M. Bennis, and M. Debbah, "Mobile unmanned aerial vehicles (UAVs) for energy-efficient Internet of Things communications," *IEEE Trans. Wireless Commun.*, vol. 16, no. 11, pp. 7574–7589, Nov. 2017, doi: [10.1109/TWC.2017.2751045](https://doi.org/10.1109/TWC.2017.2751045).
- [33] C. Pereira, J. Rodrigues, A. Pinto, P. Rocha, F. Santiago, J. Sousa, and A. Aguiar, "Smartphones as M2M gateways in smart cities IoT applications," in *Proc. 23rd Int. Conf. Telecommun. (ICT)*, Thessaloniki, Greece, May 2016, pp. 1–7, doi: [10.1109/ICT.2016.7500481](https://doi.org/10.1109/ICT.2016.7500481).
- [34] G. Wu, S. Talwar, K. Johnson, N. Himayat, and K. D. Johnson, "M2M: From mobile to embedded internet," *IEEE Commun. Mag.*, vol. 49, no. 4, pp. 36–43, Apr. 2011, doi: [10.1109/MCOM.2011.5741144](https://doi.org/10.1109/MCOM.2011.5741144).
- [35] R. Atat, L. Liu, N. Mastrorarde, and Y. Yi, "Energy harvesting-based D2D-assisted machine-type communications," *IEEE Trans. Commun.*, vol. 65, no. 3, pp. 1289–1302, Mar. 2017.
- [36] A. Afzal, S. A. R. Zaidi, D. McLernon, M. Ghogho, and A. Feki, "M2M meets D2D: Harnessing D2D interfaces for the aggregation of M2M data," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Paris, France, May 2017, pp. 1–6, doi: [10.1109/ICC.2017.7997135](https://doi.org/10.1109/ICC.2017.7997135).
- [37] G. Zhang, A. Li, K. Yang, L. Zhao, Y. Du, and D. Cheng, "Energy-efficient power and time-slot allocation for cellular-enabled machine type communications," *IEEE Commun. Lett.*, vol. 20, no. 2, pp. 368–371, Feb. 2016, doi: [10.1109/LCOMM.2015.2504980](https://doi.org/10.1109/LCOMM.2015.2504980).
- [38] S. N. Swain, R. Thakur, and S. R. M. Chebiyyam, "Coverage and rate analysis for facilitating machine-to-machine communication in LTE-A networks using device-to-device communication," *IEEE Trans. Mobile Comput.*, vol. 16, no. 11, pp. 3014–3027, Nov. 2017, doi: [10.1109/TMC.2017.2684162](https://doi.org/10.1109/TMC.2017.2684162).
- [39] S. Rao and R. Shorey, "Efficient device-to-device association and data aggregation in industrial IoT systems," in *Proc. 9th Int. Conf. Commun. Syst. Netw. (COMSNETS)*, Bengaluru, India, Jan. 2017, pp. 314–321, doi: [10.1109/COMSNETS.2017.7945392](https://doi.org/10.1109/COMSNETS.2017.7945392).
- [40] Y. Yang, G. Wu, W. Lu, and Y. Zhang, "Load balanced dynamic resource allocation for MTC relay," 2020, *arXiv:2001.01188*. [Online]. Available: <http://arxiv.org/abs/2001.01188>
- [41] S. K. Nobar, M. H. Ahmed, Y. Morgan, and S. A. Mahmoud, "Uplink resource allocation in energy harvesting cellular network with H2M/M2M coexistence," *IEEE Trans. Wireless Commun.*, vol. 19, no. 8, pp. 5101–5116, Aug. 2020, doi: [10.1109/TWC.2020.2989319](https://doi.org/10.1109/TWC.2020.2989319).
- [42] S. A. AlQahtani, "Analysis and modelling of power consumption-aware priority-based scheduling for M2M data aggregation over long-term-evolution networks," *IET Commun.*, vol. 11, no. 2, pp. 177–184, Jan. 2017, doi: [10.1049/iet-com.2016.0468](https://doi.org/10.1049/iet-com.2016.0468).
- [43] J. Huang, Y. He, Q. Duan, Q. Yang, and W. Wang, "Admission control with flow aggregation for QoS provisioning in software-defined network," in *Proc. IEEE Global Commun. Conf.*, Austin, TX, USA, Dec. 2014, pp. 1182–1186, doi: [10.1109/GLOCOM.2014.7036969](https://doi.org/10.1109/GLOCOM.2014.7036969).
- [44] W. Guo, V. Mahendran, and S. Radhakrishnan, "Achieving throughput fairness in smart grid using SDN-based flow aggregation and scheduling," in *Proc. IEEE 12th Int. Conf. Wireless Mobile Comput., Netw. Commun. (WiMob)*, New York, NY, USA, Oct. 2016, pp. 1–7, doi: [10.1109/WiMOB.2016.7763209](https://doi.org/10.1109/WiMOB.2016.7763209).
- [45] Q. T. Minh, T. K. Dang, T. Nam, and T. Kitahara, "Flow aggregation for SDN-based delay-insensitive traffic control in mobile core networks," *IET Commun.*, vol. 13, no. 8, pp. 1051–1060, 2019, doi: [10.1049/iet-com.2018.5194](https://doi.org/10.1049/iet-com.2018.5194).
- [46] Y.-B. Lin, S.-Y. Wang, C.-C. Huang, and C.-M. Wu, "The SDN approach for the aggregation/disaggregation of sensor data," *Sensors*, vol. 18, no. 7, p. 2025, 2018, doi: [10.3390/s18072025](https://doi.org/10.3390/s18072025).
- [47] P. Bosshart *et al.*, "P4: Programming protocol-independent packet processors," *Comput. Commun. Rev.*, vol. 44, no. 3, pp. 87–95, 2014, doi: [10.1145/2656877.2656890](https://doi.org/10.1145/2656877.2656890).
- [48] S.-Y. Wang, C.-M. Wu, Y.-B. Lin, and C.-C. Huang, "High-speed data-plane packet aggregation and disaggregation by P4 switches," *J. Netw. Comput. Appl.*, vol. 142, pp. 98–110, Sep. 2019, doi: [10.1016/j.jnca.2019.05.008](https://doi.org/10.1016/j.jnca.2019.05.008).
- [49] S.-Y. Wang, J.-Y. Li, and Y.-B. Lin, "Aggregating and disaggregating packets with various sizes of payload in P4 switches at 100 Gbps line rate," *J. Netw. Comput. Appl.*, vol. 165, Sep. 2020, Art. no. 102676, doi: [10.1016/j.jnca.2020.102676](https://doi.org/10.1016/j.jnca.2020.102676).
- [50] A. L. R. Madureira, F. R. C. Araújo, and L. N. Sampaio, "On supporting IoT data aggregation through programmable data planes," *Comput. Netw.*, vol. 177, Aug. 2020, Art. no. 107330, doi: [10.1016/j.comnet.2020.107330](https://doi.org/10.1016/j.comnet.2020.107330).
- [51] K. Abbas *et al.*, "An efficient SDN-based LTE-WiFi spectrum aggregation system for heterogeneous 5G networks," *Trans. Emerg. Telecommun. Technol.*, p. e3943, 2020, doi: [10.1002/ett.3943](https://doi.org/10.1002/ett.3943).
- [52] S. Anbalagan, D. Kumar, G. Raja, W. Ejaz, and A. K. Bashir, "SDN-assisted efficient LTE-WiFi aggregation in next generation IoT networks," *Future Gener. Comput. Syst.*, vol. 107, pp. 898–908, Jun. 2020, doi: [10.1016/j.future.2017.12.013](https://doi.org/10.1016/j.future.2017.12.013).
- [53] A. S. Thyagaturu, Y. Dashti, and M. Reisslein, "SDN-based smart gateways (Sm-GWs) for multi-operator small cell network management," *IEEE Trans. Netw. Service Manage.*, vol. 13, no. 4, pp. 740–753, Dec. 2016, doi: [10.1109/TNSM.2016.2605924](https://doi.org/10.1109/TNSM.2016.2605924).

- [54] A. M. Rahmani *et al.*, "Exploiting smart e-health gateways at the edge of healthcare Internet-of-Things: A fog computing approach," *Future Gener. Comput. Syst.*, vol. 78, pp. 641–658, 2018, doi: [10.1016/j.future.2017.02.014](https://doi.org/10.1016/j.future.2017.02.014).
- [55] R. Petrolo, R. Morabito, V. Loscri, and N. Mitton, "The design of the gateway for the cloud of things," *Ann. Telecommun.*, vol. 72, nos. 1–2, pp. 31–40, 2017, doi: [10.1007/s12243-016-0521-z](https://doi.org/10.1007/s12243-016-0521-z).
- [56] S. Zahoor and R. N. Mir, "Resource management in pervasive Internet of Things: A survey," *J. King Saud Univ.-Comput. Inf. Sci.*, 2018, doi: [10.1016/j.jksuci.2018.08.014](https://doi.org/10.1016/j.jksuci.2018.08.014).
- [57] M. Aazam and E.-N. Huh, "Fog computing and smart gateway based communication for cloud of things," in *Proc. Int. Conf. Future Internet Things Cloud*, Barcelona, Spain, Aug. 2014, pp. 464–470, doi: [10.1109/FiCloud.2014.83](https://doi.org/10.1109/FiCloud.2014.83).
- [58] A. U. Chowdhury and M. M. Elahi, "Design of a smart gateway for edge enabled IoT applications," in *Proc. IEEE Region 10 Symp. (TENSymp)*, Jun. 2020, pp. 417–420, doi: [10.1109/TENSymp50017.2020.9230843](https://doi.org/10.1109/TENSymp50017.2020.9230843).
- [59] *Physical Layer Aspects for Evolved Universal Terrestrial Radio Access (UTRA)*, document TS 25.814, 3GPP, Group Radio Access Network; Release 7, 2006, pp. 1–132.
- [60] T. M. de Moraes, M. D. Nisar, A. A. Gonzalez, and E. Seidel, "Resource allocation in relay enhanced LTE-advanced networks," *EURASIP J. Wireless Commun. Netw.*, vol. 2012, no. 1, pp. 1–12, Dec. 2012, doi: [10.1186/1687-1499-2012-364](https://doi.org/10.1186/1687-1499-2012-364).
- [61] M. K. Müller, F. Ademaj, T. Dittrich, A. Fastenbauer, B. R. Elbal, A. Nabavi, L. Nagel, S. Schwarz, and M. Rupp, "Flexible multi-node simulation of cellular mobile communications: The Vienna 5G system level simulator," *EURASIP J. Wireless Commun. Netw.*, vol. 2018, no. 1, p. 227, Dec. 2018, doi: [10.1186/s13638-018-1238-7](https://doi.org/10.1186/s13638-018-1238-7).
- [62] C. Mehlführer, J. C. Ikuno, M. Šimko, S. Schwarz, M. Wrulich, and M. Rupp, "The Vienna LTE simulators—enabling reproducibility in wireless communications research," *EURASIP J. Adv. Signal Process.*, vol. 2011, no. 1, p. 29, Dec. 2011, doi: [10.1186/1687-6180-2011-29](https://doi.org/10.1186/1687-6180-2011-29).
- [63] G. Rigazzi, F. Chiti, R. Fantacci, and C. Carlini, "Multi-hop D2D networking and resource management scheme for M2M communications over LTE-A systems," in *Proc. Int. Wireless Commun. Mobile Comput. Conf. (IWCMC)*, Nicosia, Cyprus, Aug. 2014, pp. 973–978, doi: [10.1109/IWCMC.2014.6906487](https://doi.org/10.1109/IWCMC.2014.6906487).



**NASSER AHMED** received the B.S. degree in computer engineering from Umm Al-Qura University, Saudi Arabia, and the M.S. degree in computer engineering from King Saud University (KSU), Saudi Arabia, where he is currently pursuing the Ph.D. degree. He is also working as a Research Assistant with KSU. His research interests include wireless networks, vertical handover in wireless networks, scheduling algorithms, resource allocation in 4G network and 5G, and M2M communication.



**NASSER-EDDINE RIKLI** (Senior Member, IEEE) received the B.S. (Ingenieur d'état) degree (Hons.) in electrical engineering (genie électrique) from INELEC, Algeria, and the M.S. and Ph.D. degrees in electrical engineering from Polytechnic University, New York. He started his academic carrier with INELEC, and then with Polytechnic University, New York, Hunter College, New York, New York City Technical College, and finally, he joined King Saud University, where he is currently working as a Full Professor with the Department of Computer Engineering. He has conducted research on funded projects with GTE Labs and Bell Labs (both in the U.S.), involving the modeling of multimedia traffic and the integration of voice and data traffic. He had many consultations with local as well as international companies and organizations in various areas of education, networking, and communication. His main research interests include performance analysis, management and administration of computer networks, both wired and wireless, design and analysis of techniques for the provision of QoS for multimedia traffic over various wired and wireless networks, queueing modeling and analysis of computer network systems, and mathematical modeling of multimedia traffic. He is the Editor-in-Chief of the *Journal of King Saud University—Computer and Information Sciences*.

• • •