

Article Multi-Connectivity for Multicast Video Streaming in Cellular Networks

Sadaf ul Zuhra^{1,*}, Prasanna Chaporkar², Abhay Karandikar^{2,†}, and H. Vincent Poor¹

- ¹ Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ 08544, USA; {sadaf.zuhra, poor}@princeton.edu
- ² Department of Electrical Engineering, Indian Institute of Technology Bombay; {chaporkar, karandi}@ee.iitb.ac.in
- * Correspondence: sadaf.zuhra@princeton.edu
- A. Karandikar is currently the Secretary to Govt. of India, Ministry of Science & Technology, Dept. of Science & Technology (on leave from IIT Bombay).

Abstract: The escalating demand for high quality video streaming poses a major challenge for communication networks today. Catering to these bandwidth hungry video streaming services 2 puts a huge burden on the limited spectral resources of communication networks, limiting the 3 resources available for other services as well. Large volumes of video traffic can lead to severe 4 network congestion particularly during live streaming events which require sending the same 5 content to a large number of users simultaneously. For such applications, multicast transmission 6 can effectively combat network congestion while meeting the demands of all the users by serving 7 groups of users requesting the same content over shared spectral resources. Streaming services can further benefit from multi-connectivity that allows users to receive content from multiple base q stations simultaneously. Integrating multi-connectivity within multicast streaming can improve the 10 system resource utilization while also providing seamless connectivity to multicast users. Towards 11 this end, this work studies the impact of using multi-connectivity (MC) alongside wireless multicast 12 for meeting the resource requirements of video streaming. Our findings show that MC substantially 13 enhances the performance of multicast streaming, particularly benefiting the cell edge users who 14 often experience poor channel conditions. We particularly consider the number of users that can 15 be simultaneously served by multi-connected multicast systems. It is observed that, about 60% of 16 the users that are left unserved under single-connectivity multicast, are successfully served within 17 the same resources by employing multi-connectivity in multicast transmissions. We prove that the 18 optimal resource allocation problem for MC multicast is NP-hard. As a solution, we present a greedy 19 approximation algorithm with an approximation factor of (1 - 1/e). Furthermore, we establish 20 that no other polynomial-time algorithm can offer a superior approximation. To generate realistic 21 video traffic patterns in our simulations, we make use of traces from actual videos. Our results 22 clearly demonstrate that multi-connectivity leads to significant enhancements in the performance of 23 multicast streaming. 24

Keywords: Multicast; Multi-connectivity; Video streaming; MBMS; 5G.

25

26

1. Introduction

Revised: Accepted: Published:

Received

1 ublished

Copyright: © 2024 by the authors. Submitted to *Network* for possible open access publication under the terms and conditions of the Creative Commons Attri- bution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/).

Citation: Zuhra, S. u.; Chaporkar, P.;

Multi-Connectivity for Multicast Video Streaming in Cellular Networks.

Karandikar, A.; Poor, H. V.

Network **2023**, *1*, 1–21. https://doi.org/

Rapid growth of video streaming applications has been the primary driver of inno-27 vation in cellular networks. As of 2023, video traffic constituted over 80% of all mobile 28 data traffic [1]. While revolutionizing the way media is consumed online, video streaming 29 has also created several challenges for telecommunication networks. Video streams are 30 resource intensive services that require a significant amount of bandwidth. As a result, the 31 exponential increase in demands for video streaming can quickly overload the network in-32 frastructure leading to network congestion which leads to slower speeds, network outages, 33 and degraded quality of service. 34

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

A large portion of video traffic is made up of live streaming from social media and 35 streaming platforms with millions of users watching the same content simultaneously. 36 These live streams pose additional challenges for the network due to their high data 37 rate, low latency, and overall quality of service requirements [2]. Using traditional one-38 to-one or unicast communications for such applications involves transmitting the same 30 content separately to each user, thus consuming a large portion of the available bandwidth. 40 Multicast transmissions are an efficient means of catering to such services by serving 41 users that need the same content simultaneously [3,4]. Multi-connectivity allows users 42 to receive content from multiple base stations simultaneously. Therefore, when a video 43 is being streamed by several base stations, allowing multi-connectivity within multicast 44 transmissions can further improve the performance of multicast streaming services. We use 45 the term *Multi-Connectivity* (MC) *multicast* to refer to such a system where multi-connectivity 46 is used alongside multicast transmissions. In this MC multicast system, users are capable of 47 multi-connectivity and can, therefore, receive multicast content from multiple base stations 48 simultaneously. 49

This paper proposes the use of MC multicast for catering to simultaneous demands for bandwidth hungry video streams. Integrating MC with multicast transmission not only boosts cell capacity but also diminishes the reliance of multicast performance on the weakest users in the system. While MC has received considerable attention for its impact on throughput and handover improvements [5–8], its unexplored integration with multicast transmissions presents a promising avenue for further research.

MC multicast allows users to potentially connect to and receive content from multiple base stations and over various Radio Access Technologies (RATs) simultaneously. It can address the demanding requirements of 5G, including high data rates, ultra-reliable low latency, and high mobility [9]. By enabling users to receive content from multiple base stations concurrently, it serves a larger user base and enhances the performance of the cell edge users. The procedures for establishing multi-connectivity within the Third Generation Partnership Project (3GPP) multicast architecture and the associated control signaling requirements have been defined in [10].

1.1. Contributions

This work studies the integration of multi-connectivity in multicast transmissions for meeting the bandwidth demands of video streaming services. We address the problem of resource allocation in a multi-connected multicast system with the aim of maximizing the number of users that can be simultaneously served using multicast transmissions. The analyses, discussions, and simulations in this work provide conclusive evidence that multi-connectivity significantly improves the performance of multicast streaming systems. The main contributions of this paper are summarized below.

- We propose a multi-connected multicast system specifically designed for video streaming. This system utilizes the existing 3GPP Multimedia Broadcast Multicast Services (MBMS) framework, enabling multicast users to receive streaming content from multiple base stations seamlessly and with minimal signaling overhead. The resulting MC multicast system serves as a low-overhead alternative to the MBMS Single Frequency Network (MBSFN) operations within 3GPP multicast systems.
- We formulate the resource allocation problem in the MC multicast system with the aim of maximizing the number of multicast users served simultaneously. Since the MC multicast system is tailored for handling concurrent demands for bandwidth-intensive video streams within limited resources, we employ the metric of *number of users simultaneously served* to measure its performance.
- We prove that the resource allocation problem in MC multicast systems is NP-hard, which means that there are no polynomial time algorithms to find the optimal solution. Therefore, we propose a centralized greedy approximation algorithm with an approximation factor of (1 - 1/e). We establish that this algorithm offers the most accurate approximation achievable for the problem.

- The centralized algorithm necessitates a central server to dictate resource allocation across all base stations within a region. Such a coordination may become impractical with increasing number of base stations. Therefore, we also propose a distributed resource allocation algorithm for MC multicast, allowing base stations to autonomously make resource allocation decisions.
- Extensive simulations clearly demonstrate the performance enhancements attained by incorporating MC in wireless multicast, particularly for video streaming applications. We employ traces from actual video streams sourced from [11] and [12] to generate realistic video traffic patterns in our simulations.

In the following section, we provide an overview of the current research across various facets of multi-connectivity and multicast in cellular mobile networks.

1.2. Related Literature

Multicast has been recognized as an effective means of catering to bandwidth hungry video transmissions [3] in cellular mobile networks. Resource allocation algorithms designed for multicast streaming have been shown to serve significantly more users while minimizing the impact of multicast streaming on other services [4,13]. Further improvements in the performance of multicast video streaming have been achieved by exploiting the inherent loss-tolerant nature of video streams [14].

The use of MC has been studied for mitigating radio link failures in ultra-dense 106 intra-frequency 5G network deployments [15], demonstrating a substantial reduction in 107 failures and throughput improvements for cell-edge users. Additionally, proportional 108 fair allocation policies have been designed [16] tailored for multi-connected ultra-dense 109 networks, prioritizing users based on load balancing and signal characteristics. MC has also 110 been shown to enhance network availability for ultra-reliable low latency communication 111 (URLLC) applications in 5G [17] where network availability is crucial. MC also optimizes 112 the system resource utilization in URLLC through load-aware cell selection [18]. 113

Numerous architectures have been proposed for implementing MC in 5G [19]. Com-114 parative evaluations in [20] assess throughput performance in distributed and cloud-based 115 heterogeneous network architectures, favoring cloud-based networks for superior through-116 put. In [21], an architecture for 5G integrating multiple RATs is proposed, facilitating 117 seamless inter-RAT MC with LTE and Wireless Local Area Network (WLAN). A control 118 and user plane split architecture for MC in 5G NR was introduced in [22], bypassing macro 119 cells for user plane transmissions of multi-connected users. It has been shown in [23] 120 that MC exhibits significant reductions in transmit power compared to single-connected 121 systems, resulting in improved outage probability and spectral efficiency. MC has also been 122 examined as a means of optimizing power consumption, particularly for 5G heterogeneous 123 cloud radio access networks [24]. Furthermore, beyond cellular networks, MC also finds 124 applications in vehicle-to-anything (V2X) services, playing a pivotal role in meeting Quality 125 of Service (QoS) requirements [25]. 126

MC combined with guard bands has also been shown to provide substantial improve-127 ment in millimeter-wave (mmWave) session continuity [26]. Methodologies in [27] evaluate 128 MC's impact on ultra-dense urban mmWave networks, showcasing enhancements in de-129 nial of service and session drop probabilities. The trade-offs between system complexity 130 and performance enhancement in multi-connected mmWave systems are explored in [28]. 131 In [29], a network throughput optimizing algorithm approaching the global optimum 132 solution was proposed for addressing the link scheduling problem in multi-connected 133 mmWave networks. Uplink MC frameworks presented in [30] efficiently monitor channel 134 dynamics and link directions in mmWave transmissions, leading to efficient scheduling 135 and session management. By mitigating radio link failures due to mobility, MC also ensures 136 seamless connectivity for mobile users [15]. Combination of MC and network coding is 137 studied in [31] to enable the transmission of high-quality video streaming services over 138 mmWave networks. 139

89

90

91

92

93

94

95

97

98

Despite the wide-ranging applications of MC, its use in multicast streaming has not yet 140 been explored in the existing literature. This work is the first to leverage MC for this crucial 141 application and establish the improvements in system performance that are achieved by 142 using MC multicast for video streaming. We also address the problem of resource allocation 143 in the proposed system. The rest of this paper is organized as follows. An overview of the 144 existing 3GPP standards for multicast and MC is provided in Section 2. This is followed by 145 a discussion of how these two techniques can be used together within the current and future 146 generations of wireless mobile networks in Section 2.1. The MC multicast system model 147 and the associated resource allocation problem are discussed in Section 3. In Section 4, we 148 prove NP-hardness of the resource allocation problem and then provide an approximation 149 algorithm for it in Section 5. We then examine the use of distributed resource allocation 150 for MC multicast in Section 6. Finally, we present the simulation results in Section 7 and 151 conclude in Section 8. 152

Notation

The set of natural numbers is denoted by \mathbb{N} . The cardinality of a set A is denoted by ¹⁵⁴ |A|. The set of integers up to n is denoted as $[n] = \{1, 2, ..., n\}$. An overview of the most ¹⁵⁵ commonly used variable notations can be found in Table 1.

 Table 1. Notation of the most commonly used variables

Symbol	Explanation
М	Number of UEs in the system
С	Number of cells/base stations in the system
N	Number of PRBs available for allocation in each cell
R	Rate of transmission of the multicast content
$r_{ik}^{c}[t]$	The maximum rate that UE k can decode on PRB j of cell c at time t
Ќ*	The MC multicast resource allocation problem

2. Multi-connectivity in MBMS

Multicast services were first standardized as part of the release 9 [32] of the 3GPP 158 standards as MBMS [33] and later as evolved-MBMS (eMBMS) [34], which is also a part of 159 the Fifth Generation (5G) New Radio (NR) [35] standards. Within MBMS, two modes of 160 multicast operation are defined, namely, Single Cell Point-To-Multipoint (SC-PTM) and 161 MBSFN. SC-PTM refers to the multicast mode where content is multicast to users within 162 a single cell. In the MBSFNs mode of operation, all the base stations within a designated 163 MBSFN area [36] transmit the same content in strict synchronization [33]. MBSFN transmis-164 sions necessitate precise synchronization between all base stations in the MBSFN area and 165 extended cyclic prefixes. This is crucial to enhance service quality for cell-edge receivers, 166 as it enables the combination of signals from various base stations, resulting in improved 167 user experience. However, the extended cyclic prefix reduces system throughput, and 168 the requirement for tight synchronization results in significant control overheads. MC 169 multicast overcomes these limitations with a considerably simpler framework than MBSFN 170 and lower transmission overheads. We discuss this in greater detail in Section 2.2. 171

The supporting architecture for MBMS with 5G NR is shown in Figure 1. The net-172 work elements that support MBMS services are the Broadcast Multicast Service Centre 173 (BM-SC), the MBMS GateWay (MBMS-GW) and the Multicell/Multicast Coordination 174 Entity (MCE) [33]. The BM-SC serves as an interface between the core network and the 175 multicast/broadcast content providers. It is responsible for transporting MBMS data into 176 the core network, managing group memberships and subscriptions and charging for MBMS 177 sessions [32]. The MCE is responsible for allocating radio resources to the base stations for 178 MBSFN operations. The MBMS-GW uses IP multicast to forward the MBMS session data 179

156 157



to the base stations. The base stations can then transmit the data to the User Equipments (UEs) via wireless multicast/broadcast.

Figure 1. MBMS architecture

In the following section, we discuss the features of MBMS that enable the use of the proposed MC multicast operations.

2.1. Enabling Multi-connectivity in Multicast Transmissions

The MBMS user plane protocol architecture defines a Synchronization (SYNC) protocol 185 layer on the transport network layer for content synchronization [37]. This layer carries 186 information needed for identifying transmission times and detecting packet loss. The SYNC 187 protocol is terminated in the BM-SC and the base stations. As a result, the MBMS content 188 sent to the base stations associated with the same BM-SC are synchronized. Consequently, 189 UEs can receive and combine multiple copies of the same content received from these base 190 stations without the need to exchange any additional control signaling. The proposed 191 multi-connectivity multicast leverages this inherent synchronization in MBMS systems, 192 enabling UEs to obtain multicast content from multiple sources without requiring additional 193 synchronization. Furthermore, since MBMS operates as an idle mode procedure, UEs can 194 use MC multicast without establishing a Radio Resource Control (RRC) connection to a 195 base station. The signaling procedures for enabling MC in MBMS have been proposed 196 in [10]. 197

For enabling multi-connectivity in MBMS, we redefine the dynamic between the primary and secondary base stations of a multi-connected UE compared to what is traditionally defined for unicast transmissions [38]. Specifically, we propose the following:

- *Connectivity:* Firstly, depending on its capability, a UE can connect to any number of base stations and receive multicast content from all of them. A UE can also remain in the RRC idle mode if it is not connected to any base station and still receive content from any number of base stations [10].
- Primary and Secondary Base Stations: For a UE using MC multicast in RRC idle mode, 205 the *primary* base station refers to the base station that it is camped on. For a UE in 206 RRC connected mode, the *primary* base station refers to the one it is connected to. All 207 other base stations from which the UEs may receive content are termed as *secondary* 208 base stations. Furthermore, the primary and secondary base stations of a UE do 209 not operate in a traditional master-slave configuration in MC multicast. Secondary 210 base stations are not dictated by the primary base station in their interaction with 211 the UE [10]. A multicast UE can receive relevant control information and multicast 212 data from multiple base stations independently. Thus, there is no real distinction 213 between primary and secondary base stations for a UE. Each base station that serves 214

the UE under MC multicast is equivalent from the perspective of the MC multicast 215 transmissions. 216

2.2. MC multicast versus MBSFN

5G NR uses MBSFN to enhance system efficiency by simultaneously transmitting 218 identical content over the same radio resources within neighboring cells grouped in an 219 MBSFN area. By leveraging the use of multi-connectivity, MC multicast can provide 220 the same advantages as MBSFN transmissions while employing a significantly simpler 221 framework with reduced transmission overheads. Similar to MBSFNs, a UE can receive 222 multicast content from multiple base stations, leading to an enhanced Signal-to-Noise 223 Ratio (SNR), particularly for cell-edge users. However, unlike MBSFN operations, base 224 stations in MC multicast are not obligated to use the same Physical Resource Blocks (PRBs) 225 for streaming multicast content. In MC multicast, identical MBMS services are streamed 226 through multiple base stations, and each base station independently allocates PRBs to the 227 multicast streams. Consequently, each base station can optimize resource allocation for 228 various services within its cell, resulting in significant frequency diversity that improves 229 the probability of reliably receiving MBMS content. A multicast UE has the flexibility to 230 decode any of the multiple copies of the content it receives. As demonstrated in Section 7, 231 this diversity leads to substantial performance improvements in terms of the number of 232 UEs served and the number of packets successfully delivered. 233

In the following section, we discuss the resource allocation problem in the MC multicast system.

3. Resource Allocation in MC Multicast

Consider a system of C cells, each with one base station serving it. There are M multi-237 connected multicast UEs in the system that can receive multicast content from any subset 238 of the *C* base stations. The set $[C] = \{1, 2, ..., C\}$ denotes the set of all cells/ base stations 239 and the set of all users is denoted by $[M] = \{1, 2, \dots, M\}$. Resource allocations decisions 240 are taken at every time slot t. In each time slot, there are N PRBs available for allocation 241 in each cell. The set of all PRBs is denoted by $[N] = \{1, 2, \dots, N\}$. We assume that there 242 is multicast content available in all the cells that is being streamed by all the UEs. The 243 multicast content is streamed at a rate of R bits per second. The UEs can potentially receive 244 the multicast streaming content from any number of neighboring base stations in addition 245 to their respective primary base stations. The multicast stream is allocated one PRB in each 246 cell, in each time slot. Resource allocation decisions are either taken independently by each 247 base station or by a central entity such as the MCE that manages the base stations within a 248 region. 249

The channel states of UEs vary as a function of time *t* as well as the PRB $j \in |N|$. On 250 PRB *j* of cell *c* at time *t*, UE *k* can decode a maximum rate of $r_{ik}^{c}[t]$ bits per second which is a 251 function of the channel state of the UE. That is, the better the channel experienced by UE k, 252 the higher will be the rate $r_{ik}^{c}[t]$. Since the multicast content is transmitted at rate R bits per 253 second, a UE may not successfully receive the multicast content from the base station that 254 it is connected to. For instance, consider that the PRB *j* is allocated to the multicast stream 255 in cell *c* at time *t*. UE *k* will be able to decode the content sent by *c* only if $R \leq r_{ik}^{c}[t]$. On the 256 other hand, if $R > r_{ik}^c[t]$, UE k will not be able to successfully decode the content sent from 257 cell c. Thus, in the absence of multi-connectivity, a UE can successfully receive data only if 258 it can decode the content from its primary cell whereas, a multi-connected UE successfully 259 receives data if it can decode the content from any one of the base stations that it receives 260 the content from. 261

Remark. Note that, even though we assume a constant bit rate *R*, video streaming traffic 262 typically uses a variable bit rate (VBR) encoding, which means that the amount of data to 263 be transmitted for the video varies over time. We employ a constant rate model for the 264 sake of simplicity in defining the resource allocation problem. However, our problem, as 265 well as the proposed resource allocation policies can be easily adapted to the VBR model 266

217

235 236

by considering the proposed setup as a snapshot of a longer VBR video stream. More specifically, to adapt to the VBR model, the transmission rate R can be made a function of time t (denote by R(t)). Then, the system model discussed above essentially represents a small enough block of time during which the rate R(t) is constant. Similarly, the resource allocation problem can be defined with the time dependent rate R(t). As we will see in the following sections, the proposed policies take allocation decisions in every time slot. Thus, the proposed policies can be used as is with the relevant rate R(t) in each time slot.

In the following, we define the resource allocation problem for this system.

3.1. Problem Definition

The resource allocation problem within the MC multicast system aims to maximize 276 the number of UEs served in each time slot. We choose the number of UEs served as 277 the optimization metric for this problem to capture the unique requirements of the MC 278 multicast problem. The primary objective of the MC multicast system is to ensure that the 279 multicast video stream is delivered to a large audience without causing network congestion. 280 Note that our system model construction ensures that only one resource is allocated to the 281 multicast stream in each time slot, which prevents overloading the system while serving 282 several video streams. Therefore, we use the number of users served to illustrate the 283 effectiveness of the resource allocation algorithms in meeting the video streaming demands 284 of users within the limited resources. 285

In the system under consideration, since a UE can receive the same content from several base stations, its performance is impacted by the resource allocation decisions across multiple cells. Therefore, the resource allocation needs to be optimized over all the C cells in the system. Throughout this paper, we assume that the users are static and do not change positions for the entire duration of the multicast transmissions. 200

For the mathematical formulation of the resource allocation problem, we first define the following sets. Assuming that every UE is trying to receive the multicast content from the base station of cell *c*, denote by $U_{jc} \subseteq [M]$ the set of users that would successfully receive the multicast content if PRB *j* is allocated to the multicast service in cell *c*, i.e., for all $c \in [C]$ and all $j \in [N]$, the set U_{jc} is given by

$$U_{ic} = \{k \in [M] : R \le r_{ik}^{c}[t]\}.$$
(1)

The collection of all such sets corresponding to cell *c* is given by

$$\mathcal{U}_{c} = \{ U_{1c}, U_{2c}, \dots, U_{Nc} \}.$$
⁽²⁾

Let \mathcal{U} be the collection of sets $\mathcal{U} = {\mathcal{U}_1, \dots, \mathcal{U}_C}$. Using these definitions, the resource allocation problem for the MC multicast system can now be stated as follows: 292

Definition 1 (Resource allocation problem \mathbf{K}^*). Given the universal set of all the users [M] and the collection of sets $\mathcal{U} = \{\mathcal{U}_1, \ldots, \mathcal{U}_C\}$, determine $\mathcal{U}' \subseteq \mathcal{U}$ such that $|\bigcup_{U_{jc} \in \mathcal{U}'} U_{jc}|$ is maximized subject to:

$$|\mathcal{U}'| = C, and \tag{3}$$

$$|\mathcal{U}' \cap \mathcal{U}_c| = 1$$
, for all $c \in [C]$. (4)

Then, in each cell $c \in [C]$ *, the PRB assigned to the multicast stream is given by* $j \in [N]$ *such that* $U_{jc} \in U'$.

The objective of the of the resource allocation problem \mathbf{K}^{\star} in Definition 1 is to maximize the cardinality of the union of sets $\bigcup_{U_{jc} \in \mathcal{U}'} U_{jc}$, which is the set of users successfully served. The solution \mathcal{U}' of \mathbf{K}^{\star} is subject to the following constraints:

275

274

- 1. For all $c \in [C]$, $|\mathcal{U}' \cap \mathcal{U}_c| = 1$: This constraint ensures that there is precisely one set 298 U_{ic} in \mathcal{U}' corresponding to each cell $c \in [C]$. That is, only one PRB is assigned to the 200 multicast stream in each cell, as required by the problem formulation. 300
- 2. $|\mathcal{U}'| = C$: This constraint ensures that there are precisely C of the U_{ic} sets in the 301 solution set \mathcal{U}' . Together with the constraint in 1, this guarantees that a set U_{ic} is 302 chosen for every cell *c*, i.e., a PRB is allocated for multicast streaming in every cell.

The resource allocation decisions are a function of: a) the channel states of UEs in each of 304 the N PRBs, b) the number of UEs streaming the multicast content, and c) the location of 305 the UEs with respect to each base station. 306

4. Computational Complexity

We show that the resource allocation problem \mathbf{K}^* is NP-hard problem and therefore, 308 no polynomial time algorithms exist for solving it. We prove this by reduction from the 309 Maximum Coverage Problem (MCP) [39], which is a known NP-hard problem defined as 310 follows. 311

Definition 2 (Maximum Coverage Problem (MCP)). Consider a universal set S, a number 312 $k \in \mathbb{N}$ and a collection of sets $\mathcal{T} = \{T_1, T_2, \dots, T_m\}$ where, for all $j \in [m], T_j \subseteq S$. The objective 313 of the MCP is to determine a sub-collection $\mathcal{T}' \subseteq \mathcal{T}$ such that $\mathcal{T}' \in \arg \max_{|\mathcal{T}'| \leq k} |\bigcup_{T_i \in \mathcal{T}'} T_i|$. 314

That is, given a collection \mathcal{T} of *m* subsets of a universal set \mathcal{S} , the objective of the MCP 315 is to find the sub-collection of at most k subsets from \mathcal{T} that cover the maximum number of 316 elements from the universal set S.

Theorem 1. The MC multicast resource allocation problem \mathbf{K}^* is NP-hard.

Proof. The proof of NP-hardness of \mathbf{K}^{\star} can be accomplished in the following steps:

- 1. First, we show that an instance of a known NP-hard problem (MCP in this case) can be reduced to an instance of \mathbf{K}^{\star} in polynomial time. This means that, we can design a polynomial time algorithm which takes MCP as the input and results in an instance of K^{\star} .
- 2. Next, we show that a solution of K^* can be mapped to a corresponding solution for the MCP in polynomial time.
- 3. Finally, using 1 and 2, we prove that no-polynomial time algorithm exists for solving \mathbf{K}^{\star} because such an algorithm would also provide a polynomial time solution for the MCP which is known to be NP-hard.

We begin by defining an algorithm to reduce an instance of MCP to an instance of K^* in polynomial time. An instance of MCP can be reduced to an instance of K^* as follows.

- Given an instance of the MCP in Definition 2 with the universal set S, the collection of • *m* sets $\mathcal{T} = \{T_1, T_2, \dots, T_m\}$ with $T_j \subseteq S$ and some $k \in \mathbb{N}$.
- Define a MC multicast system with the set of UEs [M] = S, number of cells C = k, the . number of PRBs in each cell N = m, and for all $c \in [C]$, the set $U_{ic} = T_i$.
- This defines a resource allocation problem of the form of K^* in Definition 1. This 335 reduction can be accomplished in constant time ($\mathcal{O}(C)$). 336

The pseudo-code of the algorithm for accomplishing this reduction is given in Algo-337 rithm 1. 338

307

303

317

318

319

320 321 322

323

324

325

326

327

328

329

330

331

332

333

330

343

345

346

347

348

349

350

351

352

353

354

Input: MCP with collection of sets $\mathcal{T} = \{T_1, T_2, \dots, T_m\}$ with $T_j \subseteq S$ and a number, $k \in \mathbb{N}$ **Output:** An instance of **K**^{*} with 1 $[M] \leftarrow S$ 2 $C \leftarrow k$ 3 $N \leftarrow m$ 4 for $j \leftarrow 1$ to m do for $c \leftarrow 1$ to C do 5 $U_{ic} \leftarrow T_i$ 6 7 end 8 end

This gives a one-to-one correspondence between an instance of MCP and an instance 340 of \mathbf{K}^{\star} which completes the first step of the proof. We now proceed to show that a solution 341 of the resulting instance of \mathbf{K}^{\star} can be mapped to a solution of MCP in polynomial time. 342

Let us assume that there exists a polynomial time algorithm for solving the instance of \mathbf{K}^{\star} resulting from Algorithm 1 that provides a solution \mathcal{U}' . Then, the following hold true 344 by the definition of **K***:

- $|\mathcal{U}'| = k$
- For all $c \in [k]$, $|\mathcal{U}' \cap \mathcal{U}_c| = 1$,
- \mathcal{U}' maximizes $|\bigcup_{U_{ic}\in\mathcal{U}'}U_{jc}|$.

This solution can be mapped to a solution of MCP as follows. Given the MCP in Definition 2, construct the solution set $\mathcal{T}' = \{T_1, T_2, \dots, T_m\}$ such that, if $U_{ic} \in \mathcal{U}'$, then $T_i \in \mathcal{T}'$. Since $|\mathcal{U}'| = k$, it holds that $|\mathcal{T}'| \leq k$. Therefore, by Definition 2, the constructed set \mathcal{T}' is a feasible solution of MCP. The pseudo-code for this mapping is given in Algorithm 2.

Algorithm 2: Pseudo-code for mapping a solution of K [*] to a solution of MCP		
Input: Solution of $\mathbf{K}^{\star} \mathcal{U}' \subseteq \mathcal{U}$ such that $ \mathcal{U}' = C$ and $ \mathcal{U}' \cap \mathcal{U}_c = 1$, $\forall c$		
Output: Solution of MCP \mathcal{T}'		
1 for $j \leftarrow 1$ to m do		
2 if $U_{ic} \in \mathcal{U}'$ for some c then		
$3 \mid \mathbf{T}_j \in \mathcal{T}'$		
4 end		
5 end		

To complete the proof, what is left to prove is that the constructed solution \mathcal{T}' is indeed the optimal solution of MCP. We prove this by contradiction as follows.

Let us assume that \mathcal{T}' is not the optimal solution of MCP. This implies that, there exists a set $\mathcal{T}'' \subseteq \mathcal{T}$ such that $|\mathcal{T}''| \leq k$

> $\left|\bigcup_{T_i\in\mathcal{T}''}T_j\right|>\left|\bigcup_{T_i\in\mathcal{T}'}T_j\right|.$ (5)

If (5) is true, then we can construct another solution to \mathbf{K}^{\star} , $\mathcal{U}^{\prime\prime}$ using $\mathcal{T}^{\prime\prime}$ as follows. Say $\mathcal{T}'' = \{T_{j_1}, \ldots, T_{j_\ell}\}$ with $\ell \leq k$ and say $j_1 < j_2 < \ldots < j_\ell$. We construct the set \mathcal{U}'' as follows

$$\mathcal{U}'' = \{ U_{j_1 1}, U_{j_2 2}, \dots, U_{j_{\ell} \ell}, U_{1(\ell+1)}, \dots, U_{1C} \}.$$
(6)

Then, by Definition 1, the following hold true:

- $|\mathcal{U}''| = C,$
- For all $c \in [C]$, $|\mathcal{U}'' \cap \mathcal{U}_c| = 1$, and 357
- 355 356

 $|\bigcup_{U_{ic}\in\mathcal{U}''}U_{jc}|>|\bigcup_{U_{ic}\in\mathcal{U}'}U_{jc}|,$ ٠

which contradicts our assumption that \mathcal{U}' is the optimal solution of \mathbf{K}^* . This implies that, 359 there does not exist any set \mathcal{T}'' such that $|\mathcal{T}''| \leq k$ and $|\bigcup_{T_i \in \mathcal{T}'} T_j| > |\bigcup_{T_i \in \mathcal{T}'} T_j|$. Therefore, 360 \mathcal{T}' is indeed the optimal solution of MCP. 361

Algorithm 2 maps a solution of \mathbf{K}^* to a solution of MCP in constant time ($\mathcal{O}(C)$) assignments). Thus, a polynomial time solution for \mathbf{K}^{\star} also provides a polynomial time solution for MCP. This is not possible unless P = NP. This implies that no polynomial time algorithm exists for solving \mathbf{K}^* and therefore, \mathbf{K}^* is an NP-hard problem.

Since the MC multicast resource allocation problem is NP-hard, we cannot construct a 366 polynomial time algorithm to determine its optimal solution. Therefore, in the following 367 section, we construct approximation algorithms that provide some performance guaran-368 tees. 369

5. Centralized Greedy Approximation Algorithm

We propose a greedy approximation algorithm for solving the resource allocation 371 problem \mathbf{K}^* . The Centralized Greedy Approximation (CGA) works iteratively by maxi-372 mizing the number of additional users served in each iteration. In the first iteration, CGA 373 chooses the set U_{ic} of the form in (1) from \mathcal{U} that has the largest number of elements. In the 374 subsequent steps, it picks U_{ic} 's that serves the maximum number of as yet unserved users. 375 In each step, the set chosen is from a different sub-collection \mathcal{U}_c i.e., c in the subscript of the 376 chosen sets is different for each set picked by the algorithm. The collection of sets chosen 377 after *C* iterations \mathcal{U}_{G} , is the output of the algorithm.

The steps involved in the decision making of the CGA policy are explained below. To begin, we have an empty solution set U_G .

- In the first step of CGA, the algorithm finds the largest set $U_{j^{\star}c^{\star}} \in \mathcal{U}$, i.e., $(j^{\star}, c^{\star}) \in$ 1. $\operatorname{arg\,max}_{i\in[N],c\in[C]}\{U_{jc}\}.$
- 2. The solution set \mathcal{U}_G is updated to $\mathcal{U}_G \cup \{U_{j^*c^*}\}$. This implies that the PRB j^* is allocated to the multicast stream in cell c^* .
- 3. Next, for all $j \in [N]$, the sets U_{ic^*} are removed from the set \mathcal{U} . This step ensures that the algorithm finds a feasible solution that satisfies the constraint (4) in Definition 1.
- 4. In the next step, CGA picks the set $U_{i^{\star}c^{\star}} \in \mathcal{U}$ that contains the maximum number of 387 UEs that were not present in any set U_{ic} picked in the previous iterations and assigns 388 PRB j^* to the multicast stream in cell c^* . Following this, steps 2, 3, and 4 are repeated 389 (C-1) times to determine the solution.

At the end of C iterations of CGA, the output set \mathcal{U}_G contains exactly C sets of the 391 form U_{jc} . The PRB assigned to the multicast stream in cell *c* is given by $j \in [N]$ such that 392 $U_{jc} \in \mathcal{U}_G.$ 393

The pseudo-code for this algorithm is given in Algorithm 3.

Algorithm 3: Centralized Greedy Approximation Algorithm for K [*]			
Input: Universe $[M]$, $\mathcal{U} = {\mathcal{U}_1, \ldots, \mathcal{U}_C}$, <i>C</i>			
1 Initialize: $\mathcal{U}_G = \phi$			
2 for $n = 1 : C$ do			
Pick $U_{j^{\star}c^{\star}} \in \mathcal{U}$ that covers the maximum number of elements from			
$[M] \setminus \bigcup_{U_{jc} \in \mathcal{U}_G} U_{jc}$			
$4 \qquad \mathcal{U}_G \leftarrow \mathcal{U}_G \cup \{U_{j^\star c^\star}\}$			
5 $\mathcal{U} \leftarrow \mathcal{U} \setminus \mathcal{U}_{c^{\star}}$			
6 end			

In the following theorem, we prove that the solution to K^* given by CGA has an 395 approximation factor of $\left(1-\frac{1}{e}\right)$. This means that the solution provided by this approximation 396

358

362

363

364

365

370

378 379

380

381

382

383

384

385

386

390

algorithm serves at least $\left(1-\frac{1}{e}\right)$ of the number of users that would be served by the optimal algorithm.

To state this result, we first define the following notation. Let *OPT* denote the optimal solution to the resource allocation problem \mathbf{K}^* , i.e., the optimal algorithm would serve *OPT* number of UEs in the system. Denote by m_n the number of UEs served up to the n^{th} iteration by the CGA algorithm. The gap between the optimal solution and the intermediate solution of the CGA algorithm after the n^{th} iteration is given by

$$b_n = OPT - m_n. \tag{7}$$

Therefore, $m_0 = 0$, $b_0 = OPT$ and the total number of UEs served by CGA algorithm at the end of *C* iterations is given by m_C . Using these notations, the following theorem gives the approximation factor for the CGA algorithm.

Theorem 2. The CGA algorithm (Algorithm 3) is a $\left(1-\frac{1}{e}\right)$ approximation for the resource allocation problem \mathbf{K}^* . That is,

$$m_{\rm C} \ge \left(1 - \frac{1}{e}\right) OPT.$$
 (8)

In fact, no other algorithm can achieve a better approximation unless P = NP.

To prove Theorem 2, we first prove the following two results. First, in Lemma 1, 403 we determine the lower bound on the incremental improvements in solution achieved in 404 the intermediate steps of the CGA algorithm. This result will quantify the rate at which 405 the CGA algorithm approaches the optimal solution. Then, in Lemma 2, we provide an 406 upper bound on b_n which quantifies the gap between the optimal solution *OPT* and the 407 intermediate solution of the CGA algorithm at the *n*-th iteration, m_n . Finally, using these 408 two results, we can prove that solution of the CGA algorithm is at least within $\left(1-\frac{1}{e}\right)$ of 409 the optimal solution. 410

Lemma 1. Under the CGA algorithm, the number of additional UEs served from iteration *n* to n + 1 is lower bounded by $\frac{b_n}{C}$. That is, for all $n \ge 0$, it holds that,

$$m_{n+1} - m_n \ge \frac{b_n}{C},\tag{9}$$

where C is the total number of cells in the system and b_n is in (7).

Proof. Let $U_{OPT} = \{U_1^*, \ldots, U_C^*\}$ be the optimal solution of the resource allocation problem \mathbf{K}^* where, for all $c \in [C]$, the set U_c^* is the set of UEs served by the base station of cell c. Denote by M_n , the set of users served at the end of the n^{th} iteration of CGA and $M_n^C = [M] \setminus \{M_n\}$ is the set of users not yet covered at the end of the n^{th} iteration. Then, it holds that:

$$\sum_{c=1}^{C} \left| U_{c}^{\star} \bigcap M_{n}^{C} \right| \geq \left| \bigcup_{c=1}^{C} \left(U_{c}^{\star} \bigcap M_{n}^{C} \right) \right|$$
(10)

$$\geq OPT - m_n = b_n. \tag{11}$$

Due to multi-connectivity, the sets $U_1^{\star}, \ldots, U_C^{\star}$ are not disjoint which implies the inequality ⁴¹² in (10). The quantity $\left|\bigcup_{c=1}^{C} (U_c^{\star} \cap M_n^C)\right|$ on the right hand side of (10) gives the number of ⁴¹³ unserved UEs after *n* iterations that would be served by the optimal solution U_{OPT} . To ⁴¹⁴ arrive at the inequality in (11), note that, the UEs served by the CGA algorithm may not ⁴¹⁵

402

be the same UEs that the optimal algorithm serves. Therefore, $\left|\bigcup_{c=1}^{C} (U_{c}^{\star} \cap M_{n}^{C})\right|$ is at least equal to $OPT - m_{n}$.

From (10) and (11), it follows that

$$\max_{c \in [C]} \left| U_c^* \bigcap M_n^C \right| \ge \frac{(OPT - m_n)}{C}$$
(12)

$$=\frac{b_n}{C}.$$
 (13)

Since CGA picks the set that serves the maximum possible number of yet unserved users in each iteration, we have:

$$m_{n+1} - m_n \ge \max_{c \in [C]} \left| U_c^* \bigcap M_n^C \right|. \tag{14}$$

From (13) and (14), it follows that

$$m_{n+1} - m_n \ge \frac{b_n}{C},\tag{15}$$

which completes the proof. \Box

Lemma 2. The difference between the number of UEs served in the optimal solution and the number of users served by the at the end of n + 1 iterations of the CGA algorithm is upper bounded as follows.

$$b_{n+1} = OPT - m_{n+1} \le \left(1 - \frac{1}{C}\right)^{n+1} OPT,$$
 (16)

where m_{n+1} denotes the total number of UEs served by CGA up to and including the $(n+1)^{th}$ 419 *iteration.*

Proof. We prove this result by induction. For n = 0, if

$$b_1 = OPT - m_1 \le \left(1 - \frac{1}{C}\right)OPT,\tag{17}$$

it implies that,

$$m_1 \ge \frac{OPT}{C} = \frac{b_0}{C},\tag{18}$$

which is true due to Lemma 1. Thus, the result holds for n = 0. Now we assume that

$$b_n \le \left(1 - \frac{1}{C}\right)^n OPT,\tag{19}$$

and prove the corresponding inequality for b_{n+1} . From the definition of b_n , it follows that:

$$b_{n+1} = OPT - m_{n+1} (20)$$

$$= (b_n - m_n) - m_{n+1} \tag{21}$$

$$=b_{n}-(m_{n+1}-m_{n}),$$
(22)

$$\leq b_n - \frac{b_n}{C} = b_n \left(1 - \frac{1}{C} \right), \tag{23}$$

$$\leq \left(1 - \frac{1}{C}\right)^{n+1} OPT,\tag{24}$$

418

422

where, the inequality in (23) follows due to Lemma 1 and (24) follows from (23) due 423 to (19). Therefore, by mathematical induction, the result holds for all n. This completes the 424 proof. 425

Using these results, we can now prove Theorem 2 as follows.

Proof. From Lemma 2, it follows that

$$b_C = OPT - m_C \le \left(1 - \frac{1}{C}\right)^C OPT.$$
(25)

In the limit as $C \rightarrow \infty$, from (25), it follow that

$$OPT - m_C \le \frac{OPT}{e},$$
 (26)

which implies that

$$m_C \ge \left(1 - \frac{1}{e}\right) OPT.$$
 (27)

That is, CGA provides a $\left(1-\frac{1}{e}\right)$ approximation for **K**^{*}.

To complete the proof of Theorem 2, it only remains to show that this is the best 428 possible approximation for K*. This can be easily seen using the following arguments. Let's 429 assume that there is an algorithm that could provide a better approximation for K^* . Then, 430 this algorithm would also provide a better approximation for MCP because, as we proved 431 in Theorem 1, a solution for \mathbf{K}^* can be mapped to a solution of MCP in polynomial time 432 using Algorithm 2. This is a contradiction since the greedy algorithm is known to be the 433 best possible approximation for MCP unless P = NP [40]. Therefore, no other algorithm can 434 provide a better approximation for \mathbf{K}^{\star} than the CGA algorithm. 436

This completes the proof. \Box

5.1. Comparison with optimal solutions

In this section, we evaluate the performance guarantees of the proposed CGA algo-438 rithm by comparing its solution with the optimal solution obtained for a smaller sized 439 problem of the form in Definition 1. For the purposes of this comparison, we use a 3 440 cell MC multicast system with 5 PRBs in each cell. To obtain the optimal solution, we 441 employ a brute force algorithm that works as follows. The brute force algorithm first lists 442 out all the possible PRB allocations for the 3 cell system. For instance, for a system with 443 2 cells and 2 PRBs denoted by p_1 and p_2 in each cell, the possible allocations would be 444 $(p_1, p_1), (p_1, p_2), (p_2, p_1), (p_2, p_2)$. Following this, the algorithm finds the total number of 445 UEs that would be served under each of these possible allocations. Finally, the output of 446 the algorithm is the allocation that serves the maximum number of UEs. 447

In Figures 2(a) and 2(b), we plot the number of UEs left unserved under the CGA 448 algorithm and the corresponding optimal value obtained using the brute force algorithm. 449 We refer to the plot corresponding to the brute force algorithm as 'Optimal' in the figures. 450 Figure 2(a) shows the number of UEs left unserved under the two algorithms as a function 451 of increasing number of UEs in the system. We observe that the solution of CGA matches 452 the optimal solution for up to 30 UEs in the system. As the number of UEs increases, up 453 to 3 additional UEs are left unserved while using the CGA algorithm as compared to 454 the optimal solution. Figure 2(b) shows the number of UEs left unserved under the two 455 algorithms as a function of increasing cell sizes. We observe that CGA serves just as many 456 UEs as the optimal solution for smaller cell sizes. As the cell sizes increase, one additional 457 UE is left unserved under the CGA algorithm as compared to the optimal solution. 458

426

435

437

These plots show that the CGA algorithm provides optimal solutions for smaller 459 systems. However, as the scale of the system increases, the solution provided by the CGA 460 algorithm becomes sub-optimal. 461



Figure 2. A comparison of the average number of users left unserved under CGA and the optimal resource allocation as a function of, a) increasing number of users, and b) increasing cell radii (number of users = 10).

Although the CGA algorithm provides provable approximation guarantees, it does so 462 while requiring the presence of a central controller that can make allocation decisions for 463 all base stations, based on the global view of the system. As the number of cells C increases, 464 such a centralized setup may lead to large communication overheads and increased delays. 465 In this case, a decentralized approach where base stations make allocation decisions in-466 dependently might be more feasible, albeit at the cost of losing on the performance of the 467 MC multicast streaming. In the following, we discuss the performance trade-offs between 468 the centralized and distributed allocation for MC multicast and propose a distributed 469 approximation algorithm for K^* . 470

6. Distributed Resource Allocation

In the absence of a centralized controller, allocation decisions are made by each 472 base station independently based only on the knowledge of its own cell. This type of 473 allocation does not fully reap the benefits of multi-connectivity. We illustrate this with 474 the following example. Consider a 2 cell system containing cells c_1 and c_2 . There are two 475 PRBs available for allocation in each cell. We denote these as P_1 and P_2 . Cell c_1 has four 476 users, $\{u_1, u_2, u_3, u_4\}$ and cell c_2 has two users $\{u_5, u_6\}$. All users are streaming the same 477 multicast content. Assume that user u_1 has a good channel only in P_1 and can successfully 478 receive content only on P_1 . Users u_3 , u_4 , u_5 and u_6 have a good channel only in P_2 and can, 479 therefore, successfully receive content only on P_2 . User u_2 has a good channel in both the 480 PRBs and would be served on either of them. Users u_1, u_3, u_4 are connected to both the 481 cells and can receive content from either of them. 482

Let us now look at the allocations that will be done by a distributed policy that 483 maximizes the number of users served in each cell independently. Cell c_1 considers the 484 users connected to its base station and allocates PRB P_2 to the stream because it serves the 485 maximum number of users, namely u_2 , u_3 and u_4 . Cell c_2 also optimizes independently 486 and allocates PRB P_2 to the stream to serve users u_3 , u_4 , u_5 and u_6 . Under this allocation, 487 user u_1 remains unserved even though it was multi-connected, since it could only receive 488 the content over PRB P_1 . On the other hand, users u_3 and u_4 receive content from both 489 the cells. In contrast, a centralized policy would take the users of both the cells under 490 consideration and allocate PRB P_2 to the multicast stream in c_2 and PRB P_1 to the stream in 491 c_1 and successfully serve all the users in the system. 492

Any centralized allocation policy, even if it is sub-optimal, will always do better in terms of the number of users successfully served than a policy which allocates resources in a distributed manner. A centralized policy does not necessarily mean that the policy is optimizing over the entire system. Any form of centralization that looks beyond just the individual cell will reap better performance than a completely uncoordinated allocation. In the following, we propose a distributed resource allocation algorithm for the MC multicast system that can be used even in the absence of a central controller.

6.1. Distributed Greedy Allocation

In the Distributed Greedy Allocation (DGA) policy, each base station allocates re-501 sources to the multicast streams by only optimizing over their individual cells. Although 502 allocating resources in a distributed manner will result in sub-optimal resource allocation 503 decisions as discussed above, a distributed policy allows base stations to make allocation 504 decisions independently. Therefore, such a policy can be used for enabling MC multicast 505 even in the absence of a central entity that can control all the base stations in a region. 506 Furthermore, in case of content that is highly delay sensitive, the signaling delays due 507 to the communication between the base stations and the central controller might not be 508 tolerable. For such applications, the DGA policy can be used to sacrifice optimality in favor 509 of lower delays. 510

The DGA policy solves the resource allocation problem $\mathbf{K}_{\mathbf{D}}^{\star}$ for each cell independently. The distributed resource allocation problem $\mathbf{K}_{\mathbf{D}}^{\star}$ is defined as follows. As in Section 3, $U_{jc} \subseteq [M]$ denotes the set of users that would successfully receive the multicast content if PRB *j* is allocated to the multicast service in cell *c*. Set $\mathcal{U}_{c} = \{U_{1c}, U_{2c}, \dots, U_{Nc}\}$ is the collection of such sets for cell *c*. The distributed resource allocation problem within each cell *c* can now be stated as follows:

Definition 3 (Distributed resource allocation problem $\mathbf{K}_{\mathbf{D}}^{\star}$). For all $c \in [C]$, given the collection of sets $\mathcal{U}_{c} = \{U_{1c}, U_{2c}, \dots, U_{Nc}\}$, determine $j^{\star} \in [N]$ such that $j^{\star} \in \arg \max_{i \in [N]} |U_{ic}|$.

To solve the distributed resource allocation problem $\mathbf{K}_{\mathbf{D}}^{\star}$, the DGA policy at each base station allocates a PRB to the multicast stream to maximize the number of users served by it. That is, PRB j^{\star} is assigned to the multicast stream in cell c if $j^{\star} \in \arg \max_{j} |U_{jc}|$.

The pseudo-code for this algorithm is given in Algorithm 4. The variable x_{jc} in 522 Algorithm 4 is an indicator random variable which is equal to 1 only when PRB *j* is allocated to the multicast stream in cell *c*. 524

Algorithm 4: Distributed Greedy Allocation algorithm		
Input: Sets $U_c = \{U_{1c}, \ldots, U_{Nc}\}$ for all $c \in [C]$		
1 Initialize: $x_{jc} = 0$ for every <i>j</i> , <i>c</i>		
2 for $c = 1 : C$ do		
3 Assign $j^{\star} = \arg \max_{j} U_{jc} $		
4 $x_{j^{\star}c} \leftarrow 1$		
5 end		

7. Simulations

We study the performance of the proposed MC multicast in an MBMS system con-526 sisting of seven urban macro cells [41]. A base station is located at the center of each cell 527 and UEs are distributed uniformly at random in the cells. To create 5G-specific physical 528 layer conditions, we create channels using the models recommended by 3GPP [42]. SNR 529 to rate mapping has also been done according to 3GPP specifications [42]. Other relevant 530 simulation parameters are given in Table 2. The cell edge users in the system are multi-531 connected to all the base stations in the system. In all the cells, one PRB is allocated to the 532 multicast stream in each time slot. Multi-connected users successfully receive a packet if 533

500

Parameters	Values
System bandwidth	20 MHz
Cell radius	250 m
Path loss model	$L = 128.1 + 37.6 \log 10(d)$, <i>d</i> in kilometers
Lognormal shadowing	Log Normal Fading with 10 dB standard deviation
White noise power density	-174 dBm/Hz
Noise figure	5 dB
Transmit power	46 dBm

Table 2. System Simulation Parameters [42]



Figure 3. Average number of packets successfully delivered using MC multicast as a function of increasing number of users under centralized (Algorithm 3) and distributed (Algorithm 4) resource allocation algorithms.

they can decode the content from at least one of the base stations. Other users only receive 534 the multicast content from their primary base stations.

The number of packets delivered successfully and the number of UEs successfully 536 served are used as the performance metrics in these simulations. In Figure 3, we plot the 537 average number of packets successfully received by UEs under the CGA and the DGA 538 resource allocation algorithms. One packet is transmitted in every sub-frame (1 ms) and 539 we plot the average number of packets successfully received by all the UEs in the system 540 over a period of 10 second interval (10000 packets). As expected from the discussions in 541 Section 6, we observe that the centralized policy performs better than the distributed policy. 542 However, despite its distributed nature, the packet loss under the DGA algorithm is at most 543 0.3% greater than that under the CGA algorithm. Therefore, in the absence of centralized control, the DGA algorithm can provide performance close to the centralized policy. 545

In Figures 4 to 6, we compare the performance of MC multicast with that of the 546 conventional Single-Connectivity (SC) multicast transmission. For resource allocation in 547 SC multicast, we use the DGA algorithm from Section 6 and the CGA algorithm is used in resource allocation for MC multicast. Note that, since users are connected to a single base 549 station in SC multicast, the DGA algorithm provides the optimal solution for maximizing 550 the number of users served. For the plots in Figures 4 and 5, data is transmitted at a fixed 551 rate in each sub-frame. The points in these plots are obtained by averaging over 10000 552 sub-frames. 553

Figure 4(a) illustrates the number of packets successfully delivered under MC and SC 554 multicast as the number of users increases. We observe a decline in the number of packets 555 successfully delivered as the number of UEs increases. However, the number of packets 556 successfully delivered under MC multicast is much larger than that under SC multicast. 557 Figure 4(b) plots the same metric as a function of cell radius. We observe that the number 558



Figure 4. A comparison of the average number of packets (out of 10000) successfully delivered under SC and MC multicast. Resource allocation is done using the proposed CGA algorithm (Algorithm 3) and the results are plotted as a function of, a) increasing number of users, and b) increasing cell radii.

of packets successfully delivered decreases as the cell sizes increase. This is because the path loss of the cell edge users increases as the cells become larger. The key observation here is that the performance gap between MC and SC follows an increasing trend. The relative performance of MC and SC is similar to what we observe in Figure 4a.



Figure 5. A comparison of the average number of users left unserved under SC and MC multicast. Resource allocation is done using the proposed CGA algorithm (Algorithm 3) and the results are plotted as a function of, a) increasing number of users, and b) increasing cell radii.

Figures 5(a) and 5(b) plot the average number of users left unserved in a cell per sub-frame as a function of increasing number of users and cell radius respectively. The number of users left unserved increases as the number of users and cell radius increases. The performance gap between MC and SC multicast also increases as the number of users increase. We observe that, in the absence of multi-connectivity, nearly thrice as many users are left unserved.

In Figures 6(a) and 6(b), we compare the performance of MC and SC multicast while 569 serving a real-time video stream. To generate realistic video traffic patterns, we use traces 570 of a video of Tokyo Olympics that has 133121 packets (obtained from http://trace.eas.asu. 571 edu) [11]. For these simulations, the rate of transmission varies every sub-frame, according 572 to size of the video frame being transmitted. We run the simulations for the duration of the 573 video stream (133121 sub-frames) and then average the results over the entire duration of 574 transmission. From Figure 6(a), we observe that MC multicast delivers around 8000 more 575 packets successfully than SC multicast. From Figure 6(b), we observe that 10 - 20 more 576



Figure 6. Comparison of a) the average number of packets successfully delivered (out of 133121), and b) the average number of users left unserved under SC and MC multicast while transmitting a real-time video stream. Realistic video traffic patterns generated using traces of a video of Tokyo Olympics [11].



Figure 7. Comparison of a) the average number of packets successfully delivered (out of 133121), and b) the average number of users left unserved under MC multicast and MBSFN while transmitting a real-time video stream. Realistic video traffic patterns generated using traces of a video of Tokyo Olympics [11].

UEs are left unserved under SC multicast than under MC multicast. The performance gap between the two increases as the number of UEs in the system increases.

In Figures 7(a) and 7(b), we compare the performance of MC multicast with that of 579 MBSFN transmissions. Since MBSFN requires transmitting the content over the same 580 PRB in all the cells, we choose the PRB that serves the maximum number of UEs in the 581 entire system. We use traces from a real video stream (Tokyo Olympics [11]) to generate 582 realistic video traffic patterns in these simulations as well. We observe that MC multicast 583 performs remarkably better than MBSFN. It succeeds in delivering significantly greater 584 number of packets successfully and is also able to serve many more UEs than MBSFN. 585 These results validate our claims that MC multicast can provide the benefits of MBSFNs 586 while eliminating the need for strict synchronization. In fact, as observed in Figure 7, MC 587 multicast outperforms MBSFN by large margins. 588

These simulation results clearly indicate that using multi-connectivity results in significant performance enhancements in multicast systems. The flexibility of potentially receiving content from multiple base stations results in more users being served and in reduced packet loss as well. Thus, MC multicast has tremendous potential for use in video

577

streaming services. It can help alleviate the burden on network resources while serving a larger number of users simultaneously. 594

8. Conclusions

In this paper, we propose leveraging multi-connectivity (MC) for multicast trans-596 missions and prove that it results in significant performance enhancements for multicast 597 streaming services. We address the resource allocation problem in MC multicast, aiming 598 to maximize the number of concurrently served users and prove its NP-hardness. Our 599 proposed centralized greedy approximation (CGA) algorithm for MC multicast resource 600 allocation achieves an approximation ratio of (1 - 1/e). For delay sensitive applications 601 where centralized resource allocation might become infeasible, we propose a distributed 602 greedy allocation (DGA) algorithm that enables MC multicasting without coordination 603 between base stations. We show that, despite its distributed nature, the DGA algorithm 604 results in just 0.3% more packet loss compared to the centralized policy. Using rigorous 605 simulations, we conclusively demonstrate that employing multi-connectivity in multicast 606 transmissions results in increased user coverage and reduced packet losses. Furthermore, 607 we evaluate the efficacy of our algorithms in real-time video streaming applications, utiliz-608 ing traces from authentic video streams for generating realistic traffic patterns. Performance 609 comparison of the CGA algorithm with the optimal solution obtained for a smaller problem 610 size using brute force shows that it matches the optimal solution. We also demonstrate that 611 MC multicast outperforms MBSFN, eliminating the need for strict synchronization and 612 extended cyclic prefixes. 613

9. Future Research Directions

This work provides a proof of concept for integrating MC within multicast transmis-615 sions but several practical questions remain open for further research. For instance, we 616 assume that the users are static for the entire duration of the multicast transmission. The 617 impact of user mobility on the proposed algorithms remains to be studied. Allowing for 618 mobility will imply that the sets of users served under a certain allocation keep changing as 619 a function of time. Therefore, new resource allocation algorithms need to be developed that 620 can take this into consideration. Since the problem of resource allocation in MC multicast 621 is shown to be NP-hard, machine learning based algorithms can also be developed for 622 optimizing the allocation decisions. Another interesting research direction would be to 623 consider a system where a number of different multicast streams can be simultaneously 624 broadcast in a multicast region. 625

 Author Contributions: Formal analysis, S.u.Z., P.C.; Software, S.u.Z.; Supervision, P.C., A.K., H.V.P.;
 626

 Visualization, S.u.Z.; Writing—original draft, S.u.Z; Writing—review & editing, S.u.Z., P.C., A.K.,
 627

 H.V.P.; All authors have read and agreed to the published version of the manuscript.
 628

Funding:	629
Institutional Review Board Statement: Not applicable.	
Informed Consent Statement: Not applicable.	631
Data Availability Statement: The video traces used for the simulations in Section 7 are part of the Arizona State University Video Trace Library (http://trace.eas.asu.edu/) [11].	632 633
Conflicts of Interest: The authors declare no conflict of interest.	634
Abbreviations	635
'he following abbreviations are used in this manuscript:	

595

3GPP	Third Generation Partnership Project
MBMS	Multimedia Broadcast Multicast Services
MBSFN	MBMS Single Frequency Network
MC	Multi-Connectivity
PRB	Physical Resource Block
SC	Single-Connectivity
UE	User Equipment
URLLC	ultra-reliable low latency communication

References

1.	Ericsson mobility report, 2023. [Online]. Available: https://www.ericsson.com/4ae12c/assets/local/reports-papers/mobility-	639
	report/documents/2023/ericsson-mobility-report-november-2023.pdf.	640

- Hung, Y.H.; Wang, C.Y.; Hwang, R.H. Optimizing social welfare of live video streaming services in mobile edge computing. *IEEE Transactions on Mobile Computing* 2019, 19, 922–934.
- Zuhra, S.u.; Chaporkar, P.; Karandikar, A. Efficient Grouping and Resource Allocation for Multicast Transmission in LTE. In Proceedings of the IEEE WCNC, March 2017, pp. 1–6.
- Zuhra, S.u.; Chaporkar, P.; Karandikar, A. Towards Optimal Grouping and Resource Allocation for Multicast Streaming in LTE. *IEEE Transactions on Vehicular Technology* 2019, pp. 1–1. https://doi.org/10.1109/TVT.2019.2945987.
- 5. Rosa, C.; Pedersen, K.; Wang, H.; Michaelsen, P.H.; Barbera, S.; Malkamaki, E.; Henttonen, T.; Sébire, B. Dual connectivity for LTE small cell evolution: Functionality and performance aspects. *IEEE Communications Magazine* **2016**, *54*, 137–143.
- 6. Pan, M.S.; Lin, T.M.; Chiu, C.Y.; Wang, C.Y. Downlink traffic scheduling for LTE-A small cell networks with dual connectivity enhancement. *IEEE Communications Letters* **2016**, *20*, 796–799.
- Polese, M.; Giordani, M.; Mezzavilla, M.; Rangan, S.; Zorzi, M. Improved handover through dual connectivity in 5G mmWave mobile networks. *IEEE Journal on Selected Areas in Communications* 2017, 35, 2069–2084.
- 8. Wang, H.; Rosa, C.; Pedersen, K.I. Dual connectivity for LTE-advanced heterogeneous networks. *Springer Wireless Networks* **2016**, 22, 1315–1328.
- 9. Odarchenko, R.; Aguiar, R.L.; Altman, B.; Sulema, Y. Multilink approach for the content delivery in 5G networks. In Proceedings of the International Scientific-Practical Conference Problems of Infocommunications. Science and Technology, 2018, pp. 140–144.
- 10. Karandikar, A.; Chaporkar, P.; Jha, P.K.; Zuhra, S.u. Methods and Systems for Using Multi-Connectivity for Multicast Transmissions in a Communication System. U.S. 11,368,818, Jun. 2022.
- 11. Seeling, P.; Reisslein, M. Video transport evaluation with H. 264 video traces. *IEEE Commun. Surveys Tut.* 2011, 14, 1142–1165.
- 12. Van der Auwera, G.; David, P.T.; Reisslein, M. Traffic and quality characterization of single-layer video streams encoded with the H. 264/MPEG-4 advanced video coding standard and scalable video coding extension. *IEEE Trans. Broadcast.* **2008**, pp. 698–718.
- 13. Zuhra, S.u.; Chaporkar, P.; Karandikar, A. Auction Based Resource Allocation and Pricing for Heterogeneous User Demands in eMBMS. In Proceedings of the Proc. IEEE Wireless Communications and Networking Conference (WCNC), 2019, pp. 1–6.
- 14. Zuhra, S.u.; Besser, K.L.; Chaporkar, P.; Karandikar, A.; Poor, H.V. Optimal Resource Allocation for Loss-Tolerant Multicast Video Streaming. *Entropy* **2023**, *25*, 1045.
- 15. Tesema, F.B.; Awada, A.; Viering, I.; Simsek, M.; Fettweis, G.P. Mobility modeling and performance evaluation of multiconnectivity in 5G intra-frequency networks. In Proceedings of the IEEE Globecom Workshops (GC Wkshps), 2015, pp. 1–6.
- 16. Ba, X.; Wang, Y.; Zhang, D.; Chen, Y.; Liu, Z. Effective scheduling scheme for multi-connectivity in intra-frequency 5G ultra-dense networks. In Proceedings of the IEEE International Conference on Communications Workshops, 2018, pp. 1–5.
- 17. She, C.; Chen, Z.; Yang, C.; Quek, T.Q.; Li, Y.; Vucetic, B. Improving network availability of ultra-reliable and low-latency communications with multi-connectivity. *IEEE Transactions on Communications* **2018**, *66*, 5482–5496.
- 18. Ba, X.; Wang, Y. Load-aware cell select scheme for multi-connectivity in intra-frequency 5G ultra dense network. *IEEE Communications Letters* **2019**, *23*, 354–357.
- Ravanshid, A.; Rost, P.; Michalopoulos, D.S.; Phan, V.V.; Bakker, H.; Aziz, D.; Tayade, S.; Schotten, H.D.; Wong, S.; Holland, O. Multi-connectivity functional architectures in 5G. In Proceedings of the IEEE International Conference on Communications Workshops, 2016, pp. 187–192.
- 20. Michalopoulos, D.S.; Maeder, A.; Kolehmainen, N. 5G multi-connectivity with non-ideal backhaul: Distributed vs cloud-based architecture. In Proceedings of the IEEE Globecom Workshops, 2018, pp. 1–6.
- 21. Chandrashekar, S.; Maeder, A.; Sartori, C.; Höhne, T.; Vejlgaard, B.; Chandramouli, D. 5G multi-RAT multi-connectivity architecture. In Proceedings of the IEEE International Conference on Communications Workshops, 2016, pp. 180–186.
- 22. Du, L.; Zheng, N.; Zhou, H.; Chen, J.; Yu, T.; Liu, X.; Liu, Y.; Zhao, Z.; Qian, X.; Chi, J.; et al. C/U split multi-connectivity in the next generation new radio system. In Proceedings of the IEEE Vehicular Technology Conference Spring, 2017, pp. 1–5.
- Wolf, A.; Schulz, P.; Dörpinghaus, M.; Santos Filho, J.C.S.; Fettweis, G. How reliable and capable is multi-connectivity? *IEEE Transactions on Communications* 2018, 67, 1506–1520.
- 24. Saimler, M.; Coleri, S. Multi-Connectivity Based Uplink/Downlink Decoupled Energy Efficient User Association in 5G Heterogenous CRAN. *IEEE Communications Letters* **2020**, *24*, 858–862.

638

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

- Kousaridas, A.; Zhou, C.; Martín-Sacristán, D.; Garcia-Roger, D.; Monserrat, J.F.; Roger, S. Multi-Connectivity Management for 5G V2X Communication. In Proceedings of the IEEE International Symposium on Personal, Indoor and Mobile Radio Communications, 2019, pp. 1–7.
- Kovalchukov, R.; Moltchanov, D.; Begishev, V.; Samuylov, A.; Andreev, S.; Koucheryavy, Y.; Samouylov, K. Improved Session Continuity in 5G NR with Joint Use of Multi-Connectivity and Guard Bandwidth. In Proceedings of the 2018 IEEE Global Communications Conference, 2018, pp. 1–7.
- 27. Petrov, V.; Solomitckii, D.; Samuylov, A.; Lema, M.A.; Gapeyenko, M.; Moltchanov, D.; Andreev, S.; Naumov, V.; Samouylov, K.; Dohler, M.; et al. Dynamic multi-connectivity performance in ultra-dense urban mmWave deployments. *IEEE Journal on Selected Areas in Communications* **2017**, *35*, 2038–2055.
- 28. Gapeyenko, M.; Petrov, V.; Moltchanov, D.; Akdeniz, M.R.; Andreev, S.; Himayat, N.; Koucheryavy, Y. On the degree of multiconnectivity in 5G millimeter-wave cellular urban deployments. *IEEE Transactions on Vehicular Technology* **2018**, *68*, 1973–1978.
- 29. Tatino, C.; Malanchini, I.; Pappas, N.; Yuan, D. Maximum throughput scheduling for multi-connectivity in millimeter-wave networks. In Proceedings of the 2018 International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks, 2018, pp. 1–6.
- 30. Giordani, M.; Mezzavilla, M.; Rangan, S.; Zorzi, M. An efficient uplink multi-connectivity scheme for 5G millimeter-wave control plane applications. *IEEE Transactions on Wireless Communications* **2018**, *17*, 6806–6821.
- Drago, M.; Azzino, T.; Polese, M.; Stefanović, Č.; Zorzi, M. Reliable video streaming over mmWave with multi connectivity and network coding. In Proceedings of the 2018 International Conference on Computing, Networking and Communications (ICNC), 2018, pp. 508–512.
- 32. Velde, H.; Hus, O.; Baker, M. Broadcast Operation. In *LTE-The UMTS Long Term Evolution From Theory to Practice*, 2 ed.; Chichester: John Wiley & Sons Ltd, 2011; pp. 293–305.
- 33. 3GPP TS 23.246 v.17.0.0 Rel. 17. Multimedia Broadcast/Multicast Service (MBMS); Architecture and functional description, 2022-03. [Online]. Available: https://www.3gpp.org/ftp/Specs/archive/23_series/23.246/.
- 3GPP TS 36.440 v.8 Rel. 8. Evolved Universal Terrestrial Radio Access Network (E-UTRAN); General Aspects and Principles for Interfaces Supporting Multimedia Broadcast Multicast Service (MBMS) Within E-UTRAN, 2017-03. [Online]. Available: https: //www.3gpp.org/ftp/Specs/archive/36_series/36.440/.
- 35. 3GPP TR 23.757 v.1.0.0 Rel. 17. Study on architectural enhancements for 5G multicast-broadcast services, 2020. [Online]. Available: https://portal.3gpp.org/desktopmodules/Specifications/SpecificationDetails.aspx?specificationId=3621.
- 36. Gimenez, J.J.; Carcel, J.L.; Fuentes, M.; Garro, E.; Elliott, S.; Vargas, D.; Menzel, C.; Gomez-Barquero, D. 5G new radio for terrestrial broadcast: A forward-looking approach for NR-MBMS. *IEEE Transactions on Broadcasting* **2019**, *65*, 356–368.
- 37. 3GPP TS 36.300 v.15.5.0 Rel. 15. LTE; E-UTRA and E-UTRAN; Overall description; Stage 2, 2019. [Online]. Available: https://www.3gpp.org/ftp/Specs/archive/36_series/36.300/.
- 3GPP TS 37.340 v.17.6.0 Rel. 17. 5G; NR; Multi-connectivity; Overall description, 2023. [Online]. Available: https://www.3gpp. org/ftp/Specs/archive/37_series/37.340/.
- 39. Hochbaum, D.S. Approximating Covering and Packing Problems: Set Cover, Vertex Cover, Independent Set, and Related Problems. In *Approximation Algorithms for NP-Hard Problems*; PWS Publishing Company, 1997.
- 40. Feige, U. A threshold of ln n for approximating set cover. *Journal of the ACM (JACM)* **1998**, 45, 634–652.
- 41. 3GPP TR 38.901 v.17.0.0 Rel. 17. 5G; Study on channel model for frequencies from 0.5 to 100 GHz, 2022. [Online]. Available: https://www.3gpp.org/ftp/Specs/archive/38_series/38.901/.
- 42. 3GPP TS 38.214 : Radio Access Network; NR; Physical layer procedures for data, v.17.5.0 Rel. 17, 2023.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content. 729

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723