# Beyond Complexity Limits: Machine Learning for Sidelink-Assisted mmWave Multicasting in 6G

Nadezhda Chukhno *Member, IEEE*, Olga Chukhno, *Member, IEEE*, Sara Pizzi, *Member, IEEE*,
Antonella Molinaro *Senior Member, IEEE*, Antonio Iera, *Senior Member, IEEE*,
and Giuseppe Araniti, *Senior Member, IEEE*

*Abstract*—The latest technological developments have fueled revolutionary changes and improvements in wireless communication systems. Among them, mmWave spectrum exploitation stands out for its ability to deliver ultra-high data rates. However, its full adoption beyond fifth generation multicast systems (5G+/6G) remains hampered, mainly due to mobility robustness issues. In this work, we propose a solution to address the problem of efficient sidelink-assisted multicasting in mobile multimode systems, specifically by considering the possibility of jointly utilizing sidelink/device-to-device (D2D), unicast, and multicast transmissions to improve service delivery. To overcome the complexity problem in finding the optimal solution for user-mode binding, we introduce a pre-optimization step called *multicast group formation (MGF)*. Through a clustering technique based on unsupervised machine learning, *MGF* allows to reduce the complexity of solving the *sidelink-assisted multiple modes mmWave (SA3M)* problem. A detailed analysis of the impact of various system parameters on performance is conducted, and numerical evidence of the complexity/performance trade-off and its dependence on mobility patterns and user distribution is provided. Particularly, our proposed solution achieves a network throughput improvement of up to 32% over state-of-the-art schemes while ensuring the lowest computational time. Finally, the results demonstrate that an effective balance between power consumption and latency can be achieved through appropriate adjustments of transmit power and bandwidth.

*Index Terms*—6G, millimeter wave, multicast, unicast, sidelink, radio resource management, mobility, machine learning.

## I. INTRODUCTION

**F**UTURE wireless networks are anticipated to deliver a wide range of services requiring improved performance

compared to the fifth generation (5G) in terms of delivered data rate, tolerated latency, mobility support, and massive access. Such services make use of different kinds of wearable devices, including head-mounted displays, motion-tracked controllers, haptic gloves, and body-tracking sensors.

Various future mobile applications, such as camera-assisted automotive driving, virtual reality with rich sensory information, and holographic communications, call for extra-high-demanding service delivery requirements that current communication technologies, operating in the low- and middle-frequency bands, are unable to meet. Millimeter-wave (mmWave) communication is considered a viable way to overcome this challenge, enabling multi-Gigabit/s data rates and ultra-low latency for a high number of devices thanks to its wide bandwidth and the compact antenna size allowed by short wavelength communications [1].

However, relying solely on mmWave technology may not be sufficient to ensure efficient spectrum utilization for future group-oriented applications. Exploiting multicasting becomes crucial to enhancing the capabilities of mmWave communications and boosting network utility while saving spectrum resources since it enables the simultaneous delivery of content to multiple users through a single transmission, thereby optimizing network efficiency. Coupling mmWave technology with multicast transmission is becoming an important research trend toward increasing energy efficiency and network throughput [2], [3]. However, in the case of sparse user deployments, serving the whole group with either one wide beam or a set of directional beams may significantly reduce the benefits of exploiting extremely high frequencies.

Device-to-device (D2D) is a technology that can be effectively leveraged in such a scenario. Primarily, D2D involves two devices in close proximity communicating directly without a base station (BS). In 5G this is done over the sidelink, which is defined as the interface between user equipment devices (UEs) for direct communications. Thanks to the ability of sidelink-assisted communications to achieve ultra-low latency connectivity, high data rates, and ultra-high reliability [4], [5], this technology is expected to play the same fundamental role in 6G as it actually does in 5G. The main reason why D2D can assist in highly directional multicast communications is that beam narrowing can be achieved by excluding scattered users from the multicast transmission and replacing the transmission from the BS with the creation of D2D links between interacting nodes in the local range.

Generally, existing research works only focus on unicasting/multicasting in mmWave networks or provide heuristic solutions to the complex multicasting problem (sometimes enhanced by D2D) in directional systems (see, e.g., [6]) under no mobility constraint. This raises a fundamental question:

*Is there a framework for sidelink-assisted mmWave multicasting, which works for mobile scenarios and is scalable when working with a large number of users?*

In this work, we focus on multicast scheduling assisted by sidelink and unicast transmissions. In particular, we take a cue from our previous investigations on multicast data transmission optimization [7] and consider scenarios that include multicast users moving at low speeds, such as pedestrians equipped with wearable devices. Differently from [7], where we introduced a framework for mmWave beam coverage estimation, we propose to split users into groups by leveraging fast algorithms (e.g., unsupervised hierarchical clustering [8], which has received much attention in the literature in this field), and then enable the system to modify the transmission mode of moving users for performance improvement by considering their dynamic channel conditions. Considered options are *(i)* unicasting, *(ii)* sidelink unicasting/multicasting, and *(iii)* mmWave multicasting, and, for each user's group and each occurring condition, the designed policy makes the decision whether to widen the beam or use different beams.

The main contributions of our study are as follows:

- *Multicast services delivery framework:* We propose a framework that integrates sidelink/D2D, unicast transmissions, and multicasting, tailored to finding the optimal user-mode association solution. The framework serves as a mobility management tool for dynamic directional multicast systems and includes two steps: *(i) multicast group formation (MGF)* for forming multicast groups and *(ii) sidelink-assisted multiple modes mmWave (SA3M)* step for establishing D2D and unicast links.

- *Complexity reduction via unsupervised machine learning:* We leverage unsupervised machine learning techniques to cluster users in the MGF step, with the aim of reducing the complexity of finding the optimal SA3M solution.

- *Optimization problem formulation:* We formulate the transmission scheduling problem in the SA3M step as an optimization problem that maximizes network throughput. This formulation enables the development of an optimal scheduling solution.

- *Low-complexity heuristic solution:* We present a low-complexity heuristic solution that yields comparable results to the proposed scheduling solution. This approach reduces computational complexity while maintaining effective performance.

- *Analysis of complexity vs. performance trade-off:* We conduct an extensive numerical analysis to explore the trade-off between complexity and performance, considering user mobility. The results provide practical insights for achieving the desired trade-off.

- *Guidelines for transmit power tuning to reduce power consumption:* We provide numerical results illustrating the potential of adjusting transmit power and transmission bandwidth to effectively reduce total power consumption in the system.

The remainder of this paper is organized as follows. Related works are surveyed in Section II. In Section III, we explain the motivation behind the proposed two-step approach. Section IV details the system model. The proposed framework for multicast scheduling assisted by sidelink and unicast transmissions is formulated in Section V, where we also introduce a heuristic algorithm. Numerical results and algorithms' performance comparison are discussed in Section VI. Section VII concludes this work.

## II. BACKGROUND AND RELATED STUDIES

The use of directional beams, as opposed to conventional omnidirectional systems, has a significant impact on various aspects of wireless system design. Therefore, in the following, we first provide a background on multicasting with sidelink assistance to deal with mmWave propagation challenges. Then, we review sidelink-assisted multicast communications and machine learning-aided group-based mmWave transmission approaches.

### A. mmWave D2D-Aided Multicasting

D2D communications can significantly improve the coverage of systems operating in the mmWave band by serving UEs that are not covered by the directional beam via proximity communications. The problem of D2D-assisted multicasting in mmWave directional systems has been the focus of several recent studies. In [9], an efficient heuristic for multicast data delivery is developed, where multi-hop and simultaneous D2D transmissions (also known as spatial reuse) are combined to achieve reduced power consumption compared to mmWave multicasting performed through a series of unicast transmissions. In [10], the optimal multicast scheduling problem is addressed by exploiting D2D and simultaneous transmissions, multicast group partitioning, and beam selection by exploiting a multilevel codebook structure. Furthermore, in [11], D2D communication has been shown to increase the efficiency of multicasting. The authors propose a user clustering and multicast path planning algorithm with cubic complexity on the set of multicast users.

In [12], a similar approach is proposed, wherein multicast scheduling takes advantage of the relaying and spatial sharing features of mmWave networks operating at 73 GHz. The proposed multicast method reduces the total data delivery time for all multicast group members by properly selecting transmitting nodes and their target destinations at each time slot. Similarly, in [13], an optimal D2D-enabled multicast scheduling policy is proposed to minimize energy consumption in mmWave cellular networks. The authors solve the joint problem of D2D communications and concurrent transmissions for multicast data delivery, where multicast transmission from BS to the users is implemented through multi-hop D2D links. In [14], an optimal sidelink-aided multicast system for multiquality tiled 360 VR video is proposed to achieve high user experience under the constraints of bandwidth

resource and tile quality smoothness in overloaded situations. A suboptimal solution of low complexity is also offered.

Differently, the security aspect is considered in [4], wherein the proposed D2D protocol manages the efficient and reliable delivery of multicast data to a group of IoT devices. In [15], a sidelink-assisted cooperative retransmission scheme is proposed, according to which the neighboring UEs assist an error-prone downlink transmission.

A common feature of all the studies mentioned above is the fact that *all assume only a static scenario* i.e., the mobility of multicast groups is not considered, which means that there are no significant environmental changes affecting the channels between the BS and the group.

A further challenge emerges in directional multicast systems in the case of a *non-static scenario*. In mmWave multicast systems, beams are steered between users to cover multiple receivers at once, which leads either to signal degradation or even to the interruption of the connection between the BS and the mobile receiver if the latter is located near the edge of the beam. Therefore, ensuring coverage in the presence of mobile users is becoming more and more challenging. Mobility aspects in mmWave multicast systems have been considered in [16], where, based on the training information and starting only with the finest beams, a scalable beam grouping algorithm (without D2D capabilities) is designed to achieve minimum multicast group data transmission time. Then, in [17], a mode selection algorithm (cellular communication based on uplink/ downlink, including multicasting and direct vehicle-to-everything (V2X) communications using sideline), is proposed based on a heuristic with the goal of avoiding the overload in the sidelink resources. The algorithm is based on controlling the received signal on the user side.

While research efforts have focused on multi-beam systems, *this paper explicitly examines single-beam multicast transmission*. We make this choice based on the assumption that mobile devices typically have a single radio frequency (RF) chain to ensure cost efficiency [18].

### B. Machine Learning for mmWave Transmission

Machine learning (ML) techniques provide fast and efficient solutions to multicast group formation, D2D clustering, and transmission mode selection. In [8], a Self-Organizing Map (SOM) is used to perform multicast group formation, whereas the D2D technology is exploited to deal with blockages. In [19], an intelligent mode selection strategy in D2D-assisted 5G heterogeneous networks is proposed to improve the performance of virtual reality (VR) data broadcasting in terms of network throughput. The policy consists of two main parts: *(i)* fast D2D clustering algorithm based on unsupervised learning, *(ii)* smart mode selection based on reinforcement learning (Nash-Q-learning and WoLF-PHC) to find the optimal transmission strategy in every time slot among broadcasting, mmWave unicasting, and D2D multicasting.

A similar strategy is utilized in [20], where different supervised ML algorithms are executed to define the subset of users that shall be served directly by the eNB[1] instead of

---

[1] Evolved Node Bases (eNBs) is the LTE term for a BS.

D2D clusters (by default). Here, ML is an effective tool to address the identified problem since it exploits offline training without involving the eNB, thereby distributing the training process. Data (distance and channel conditions) required to create clusters in an online implementation can be obtained via the D2D discovery process, which occurs before the communication begins.

ML techniques have also been exploited to determine beam direction, beam weights, transmit power, and blockage predictions for directional unicast systems [21], [22], [23].

## III. MOTIVATIONS

In 5G systems and beyond, the user-mode association problem poses a significant challenge due to the complexity of calculating unicast, multicast, and sidelink communication mode combinations for all users.

Such complexity can be managed with the help of supervised, unsupervised, and reinforcement learning (RL) as well as optimization techniques. RL can outperform the optimization in scenarios with rapidly changing channels, coverage, and topology, i.e., a large state space [24], [25]. In this case, the problem can be formulated as multi-agent reinforcement learning (MARL) with centralized rewards, which requires addressing privacy concerns [26], [27]. For this purpose, the protocols used by users to communicate their actions to the central system must be carefully designed. In [19], the authors use MARL for a static scenario and do not consider multicasting. Since we consider mobility, the user's action is his changing position while moving. In the multicasting scenario under analysis, wherein multiple users are served with the same beam, the complexity further increases because users' actions are dependent on each other. In fact, the data rate of the group is limited by the user with the worst channel conditions, which dynamically changes due to mobility. Also, in our scenario, users do not communicate with each other as they move around and are unaware of other users' actions (for example, changes in location).

As for supervised machine learning, offline learning is not suitable due to its "configuration specific" nature. Variations of transmission power, number of users, area of interest, and other transmission parameters affect the final result. This means that we have to provide offline training for each configuration. Furthermore, optimization is still needed for the collection of datasets. Similarly, in the case of online supervised machine learning, the model also needs the training dataset and should be retrained on ground truth values over time. For the considerations discussed above, both MARL and offline/online supervised ML are not viable for mobile mmWave multicasting.

While developing an optimal strategy for multicasting in [28], we encountered complexity issues starting with 15 users in terms of computational time and RAM on disk needed to create all possible options for exhaustive search. Therefore, a *trade-off between complexity and optimality* should be considered to address scalability issues in cellular systems working in real time.

Fig. 1. System illustration.

TABLE I
COMPUTATIONAL TIME, SECONDS

| Number of users | 2 | 5 | 7 | 10 | 12 | 15 | 17 | 20 | Complexity |
|---|---|---|---|---|---|---|---|---|---|
| Optimal multicasting [28] | 0.48 | 0.6 | 3.6 | 601.8 | 3261 | 3600* | - | - | $O(2^N)$ |
| Optimal user-mode association | 1.8 | 10.8 | 1805* | - | - | - | - | - | $O((2^N)^3)$ |
| Heuristic multicasting [29], [30] | 0.007 | 0.015 | 0.017 | 0.019 | 0.024 | 0.028 | 0.032 | 0.036 | $O(N \cdot |\Theta|)$** |
| Proposed MGF | 0.0004 | 0.0004 | 0.0004 | 0.0004 | 0.0004 | 0.0004 | 0.0004 | 0.0004 | $O(N^2)$*** |
| *Suboptimal solution due to the time restrictions. | | | | | | | | | |
| **$|\Theta|$ is the number of beamwidths, such as $|\Theta| = [10° : 10° : 180°]$. | | | | | | | | | |
| ***Time complexities are general approximations and can vary based on factors such as the specific algorithm used, the distance metric employed, the efficiency of the implementation, and locations of data points. | | | | | | | | | |

Table I collects the results in terms of computational time required for: *(i) optimal multicasting* in [28], *(ii) optimal user-mode association* considering multicast, unicast, and D2D, *(iii)* the *heuristic multicasting* solution in [30], and *(iv)* the *proposed MGF* via unsupervised learning that works as a first step of the designed framework. Optimal multicasting has considerable complexity that is further compounded when unicast and D2D modes are also considered. A significant reduction in complexity is achievable when using a heuristic solution. Significant benefits can be gained by leveraging unsupervised machine learning for multicasting, as in the proposed MGF step. To implement a practical user-mode association strategy, the design choice in this work is *to leverage unsupervised ML for multicast user clustering* before system performance is optimized, considering the possibility of establishing unicast and D2D communications, as will be discussed in the following sections.

## IV. SYSTEM MODEL

This section introduces the reference scenario and describes traffic, antenna, propagation, blockage, and mobility models. The reference system is depicted in Fig. 1, while the system modeling notation is reported in Table II.

### A. Deployment and Traffic Model

We examine a 5G NR outdoor deployment, wherein all UE devices, such as XR glasses and wearable headsets, are provisioned with mmWave modules to be served by an NR BS

that operates in the 28 GHz band. The height of the NR BS is set to $h_A$, and its coverage radius is $R_d$, within which all UEs can successfully receive data. The geometric locations of UEs are assumed to be scattered across a plane. In our system, all UEs, $\mathcal{N} = \{1, \ldots, N\}$, are assumed to be dynamic.

We assume that all UEs from $\mathcal{N}$, located and moving within a specific area of interest, require the same multicast service. In practical deployments, both multicast and unicast sessions may coexist. In this work, we do not consider unicast sessions and focus on a single multicast transmission, mainly to analyze the performance of the proposed framework in case of no "external disturbances" in the system. The problem of the joint management of unicast and multicast traffic is, by itself, a research problem that deserves particular attention [31], [32].

We specify that the concepts of session type and delivery methods are different. In this work, only a multicast session (i.e., one-to-many content/data stream delivery) is considered, while both multicast and unicast transmission modes can serve multicast UEs of a given session.

### B. Antenna Model

We assume that devices transmit directionally with an antenna pattern akin to a conical shape, i.e., beamwidths are symmetric in both the vertical and horizontal planes. To this end, we approximate the beamforming pattern with the following transmit antenna gain as proposed in [33]:

$$G_{\text{tx}} = D_0 \rho(\alpha_i), \tag{1}$$

TABLE II
SYSTEM MODELING NOTATION

| Parameter | Definition |
|---|---|
| $f_c$ | Carrier frequency |
| $W$ | Available bandwidth |
| $R_d$ | Radius of area of interest |
| $h_A$ | Height of NR BS |
| $h_U$ | Height of UE |
| $h_B$ | Height of blocker |
| $r_B$ | Radius of blocker |
| $\lambda_B$ | Density of blockers |
| $N_U$ | UE planar antenna elements |
| $P_T$ | Transmit power |
| $P_{T,d}$ | D2D transmit power |
| $N_0$ | Power spectral density of noise |
| $v$ | UE's velocity (pedestrian/segway) |
| $S_{\text{thr}}$ | SINR threshold |
| $S_{\text{thr,h}}$ | SINR threshold for heuristic |
| $N$ | Number of multicast UEs |
| $M_{S,nB/B}$ | Fading margins |
| $B$ | Packet size |
| $y_i$ | Distance between UE $i$ and NR BS |
| $p_B(y_i)$ | Blockage probability of UE $i$ |
| $\mathcal{N}$ | Set of multicast UEs |
| $D_0$ | Antenna directivity |
| $\alpha_i$ | Angular deviation from UE $i$ antenna boresight |
| $G_{\text{tx}}, G_{\text{rx},i}$ | Antenna array gains at NR BS and UE $i$ |
| $\theta$ | Half-power beamwidth |
| $L_{\text{dB}}$ | Path loss in linear and decibel scales |
| $A, \varsigma$ | Propagation coefficients |
| $S$ | Signal-to-noise ratio |
| $\mathcal{G}$ | Set of all multicast groups |
| $\mathcal{G}_m$ | Set of UEs in a multicast group covered by the same beam $m$ |
| $\mathcal{U}$ | Set of unicast UEs |
| $\mathcal{L}$ | Set of all unicast and multicast groups |
| $\mathcal{G}^*$ | Set of UEs served via multicast |

where $D_0$ is the maximum antenna directivity along the antenna boresight, $\alpha_i$ is the angular deviation of the transmit/receive direction from the boresight of a directional antenna for receiver $i$, $i \in \mathcal{N}$, and $\rho(\alpha_i) \in [0;1]$ is a piecewise-defined linear function that scales the directivity $D_0$ with respect to the angular deviation [33].

### C. Propagation and Blockage Model

Following 3GPP standard [34], we exploit the 3GPP urban microcell (UMi) street canyon path-loss model:

$$L_{\text{dB}} = 32.4 + 21 \log_{10} y_i + 20 \log_{10} f_c, \qquad (2)$$

where $f_c$ is the carrier frequency in GHz, and $y_i$ is the three-dimensional (3D) distance between the BS and the UE $i$.

5G NR systems operating in a high-frequency band suffer from moving obstacles (called "blockers"), including humans and vehicles. Here, pedestrians are assumed to temporarily block the line-of-sight (LoS) path between the UE and the NR BS, i.e., causing blockage by the human body. This blockage attenuation is considered to be 15 dB. We also introduce shadow fading margins represented by $M_{S,B}$ and $M_{S,nB}$ for the blocked and non-blocked states, respectively. Then, the path loss in (2) may be written in a linear scale using $Ay_i^\varsigma$, with $A$ and $\varsigma$ being propagation coefficients:

$$A_{\text{LoS,nB}} = 10^{2 \log_{10} f + 3.24} M_{S,nB}, \; \varsigma_{\text{LoS}} = 2.1,$$

$$A_{\text{LoS,B}} = 10^{2 \log_{10} f + 4.74} M_{S,B}, \; \varsigma_{\text{LoS}} = 2.1. \qquad (3)$$

The blockers are modeled as cylinders with height $h_B$ and radius $r_B$ [35]. The number of blockers follows a Poisson distribution with density $\lambda_B$ per square meter.

Then, the signal-to-noise ratio (SNR) in the propagation model can be represented as

$$S = \frac{P_T D_0 \rho(\alpha_i)}{N_0 W} \left( \frac{y_i^{-\varsigma_{\text{LoS}}}}{A_{\text{LoS,nB}}} \big[ 1 - p_B(y_i) \big] + \frac{y_i^{-\varsigma_{\text{LoS}}}}{A_{\text{LoS,B}}} p_B(y_i) \right), \quad (4)$$

where $p_B(y_i)$ is the blockage probability at the 3D distance $y_i$, $N_0$ is the noise power spectral density, and $W$ is the operating bandwidth.

### D. Mobility Model

We assume that UEs follow the social force-based mobility model that captures the realism of crowd behaviors. More specifically, we apply the Headed Social Force Model (HSFM) proposed in [36], which can reproduce pedestrians moving together. The HSFM model allows us to test the real-life scenario composed of several groups of moving UEs (with speed $v$) and is relevant to our system as we consider the multicast content delivery for a set of UEs.

## V. PROPOSED SIDELINK-ASSISTED MULTICAST SCHEDULING

### A. Framework Description at a Glance

The primary objective of our proposed framework is to maximize the system throughput for delivering a multicast session to multiple mobile UEs by dynamically selecting the transmission mode (i.e., unicast, multicast, or sidelink) for each member of the multicast group. This section details our proposed solution for the dynamic sidelink-assisted mmWave scheduling problem. The framework consists of two steps: the *multicast group formation (MGF)* and the *sidelink-assisted multiple modes mmWave (SA3M)* scheduling. The general flow diagram relevant to our proposal is presented in Fig. 2.

First, we perform the MGF step by applying a hierarchical clustering algorithm, which is an unsupervised machine learning technique that aims to find natural grouping based on the characteristics of the input data. We accomplish this task based on the information about the UE locations. Note that different clustering methods can be used at this stage. Among them, we chose hierarchical clustering because it is characterized by low complexity and considers UEs' position, which is essential for directional multicast transmissions.

After initial clustering and multicast group configuration by MGF, as UEs move at speed $v$, the system may implement the SA3M step. The output is a change in the transmission mode [2] involving the use of *(i)* sidelink transmission(s), *(ii)* or unicast transmisison(s), or *(iii)* the beamwidth adjustment for the multicast group. Specifically, starting from the configuration

---

[2]Transmission mode switching means a change in the initial scheduled multicast transmissions for groups of multicast UEs (after MGF) so as to adapt the system to UE's mobility. For example, if a UE moves far away from its original multicast group, unicast transmission may be employed specifically for this UE, while maintaining multicast transmission for the remaining UEs.
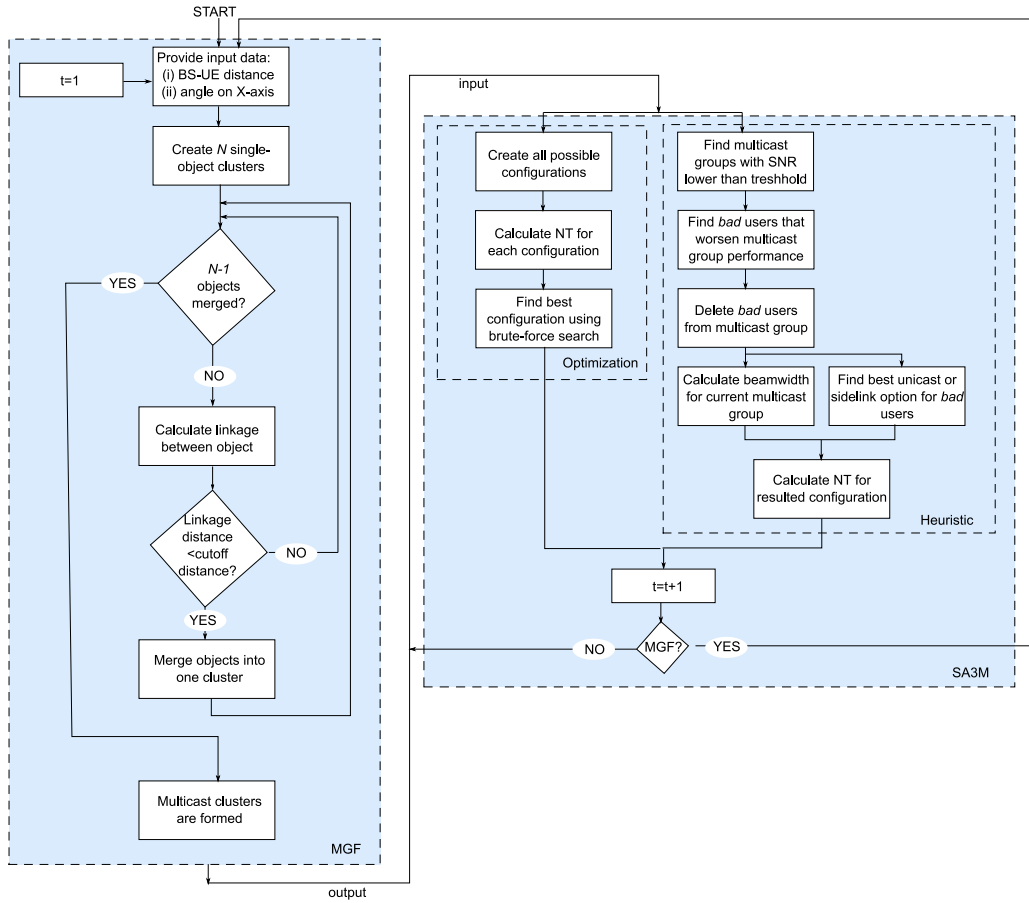
Fig. 2.  Flow diagram of proposal.

generated by MGF, the optimal SA3M algorithm exhaustively searches for all possible switching options, that is $1, 2, \ldots$, or all $N$ UEs can be served via unicast while the rest of UEs belong to multicast groups. Similarly, $1, 2, \ldots$, or $N-1$ UEs can receive the data through sidelink ($N-1$ since at least one UE has to be a relay node towards the sidelink receivers). Note that, in case a UE needs to join a different multicast cluster, MGF, which serves as the means to build multicast clusters, is rerun. Without MGF, the algorithm would need to examine all potential multicast clusters, i.e., $2^N - 1$ possible multicast groups (instead of the configuration containing just $n$ groups selected by MGF). Moreover, the procedure described above would have to be run for sidelink and unicast options for each multicast configuration, significantly complicating the model. Thus, the *main mission of MGF is to reduce the complexity of the SA3M step.*

### B. Step I – Multicast Group Formation

Hierarchical clustering, applied for MGF, builds a binary merge tree. It starts from the data elements stored at the leaves (interpreted as singleton sets) and proceeds by merging two by two the *closest* subsets (stored at nodes) until the root of the tree contains all the elements of $X$. Specifically, in the beginning, each data point is assumed to be a separate cluster. Then, similar clusters are iteratively combined. We denote by $\Delta(X_k, X_j)$ the distance between any two subsets of

$X$, called the linkage distance. This technique is also called *agglomerative hierarchical clustering* [37]. In our case, $X$ represents an array with the observations, with at least one column and $N$ strings (each string corresponds to a UE). The reference angle from the $X$-axis and the distance between the BS and every UE are used as the observations.

Let $D(x_k, x_j)$ denote the elementary distance between any two elements of $X$ (e.g., Euclidean, Minkowski, Chebyshes, etc.). In order to select the closest pair of subsets at each stage of the hierarchical clustering, we define a subset distance $\Delta(X_k, X_j)$ between any two subsets of elements. When both subsets are singletons $X_k = x_k$ and $X_j = x_j$, then $\Delta(X_k, X_j) = D(x_k, x_j)$. There are four different methods to measure the similarity between clusters, i.e., four *common* linkage functions (also known as cluster-level scoring functions) that calculate the distance between clusters:

• Single linkage (SL) represents the shortest distance among all data points in two clusters, i.e.,

$$\Delta(X_k, X_j) = \min_{x_k \in X_k, x_j \in X_j} D(x_k, x_j).$$

• Complete linkage (CL) represents the farthest distance among all data points in two clusters, i.e.,

$$\Delta(X_k, X_j) = \max_{x_k \in X_k, x_j \in X_j} D(x_k, x_j).$$

---

**Algorithm 1:** MGF

1 **Input:** $X$;
2 **Output:** Multicast clusters;
3 **Initialize** $\mathcal{G}_k = \{x_k\}$, $k = 1, \ldots, N$, $\mathcal{L} = \{\{x_1\}, \{x_2\}, \ldots, \{x_N\}\}$, *distance threshold*;
4 *counter* $\leftarrow N$;
5 **while** *counter* $\neq 2$ **do**
6      **Select** $\mathcal{G}_k$ and $\mathcal{G}_j$ from $\mathcal{L}$ such as $\Delta(X_k, X_j)$ is minimized along all pairs;
7      **if** $\Delta(X_k, X_j) <$ *distance threshold* **then**
8          **Merge** $\mathcal{G}_k \cup \mathcal{G}_j$;
9          $\mathcal{L} \leftarrow \mathcal{L} \setminus \mathcal{G}_k$;
10          $\mathcal{L} \leftarrow \mathcal{L} \setminus \mathcal{G}_j$;
11          $\mathcal{L} \leftarrow \mathcal{L} \cup (\mathcal{G}_k \cup \mathcal{G}_j)$;
12          *counter* $\leftarrow N - 1$;
13      **else**
14          go to line 5;
15 **return** $\mathcal{L}$;
16 **end**

- Average linkage uses the average distance between all pairs of objects in any two clusters, i.e.,

$$\Delta(X_k, X_j) = \frac{1}{|X_k|} \frac{1}{|X_j|} \sum_{x_k \in X_k} \sum_{x_j \in X_j} D(x_k, x_j).$$

- Ward linkage (appropriate for Euclidean distances only) uses inner squared distance, i.e.,

$$\Delta(X_k, X_j) = \sqrt{\frac{2|X_k||X_j|}{|X_k| + |X_j|}} \|\bar{x}_k - \bar{x}_j\|_2,$$

where $\|\|_2$ is the Euclidean distance, $\bar{x}_k$, $\bar{x}_j$ are the centroids of clusters $X_k$ and $X_j$, respectively.

We note that, in the case of hierarchical clustering, the number of clusters may not be determined in advance as, for example, in the case of the *k*-means algorithm. Here, either a cutoff distance or a maximum number of clusters must be specified. In this work, we exploit a cutoff distance, which is the linkage distance threshold above which clusters will not be merged. The pseudo-code of the hierarchical clustering adapted for MGF is presented in Algorithm 1.

The algorithm assigns each observation in $X$ to a single-object cluster (line 3). Then, the algorithm computes similarity information between every pair of objects $\mathcal{G}_k$ and $\mathcal{G}_j$ in the data set and uses a linkage function to group objects into a hierarchical cluster tree (line 6). Therefore, objects/clusters in close proximity are linked together if the result of the linkage function does not exceed the cutoff distance *distance threshold* (lines 7-14). This determines where to cut the hierarchical tree into clusters, thereby partitioning the data.

Since we start from *counter* $= |X| = N$ and finish with a root containing the full set $X$, the algorithm performs exactly $N - 1$ merge operations. A straightforward implementation yields a cubic time complexity, in $O(N^3)$, since, in the *k*-th iteration of $N - 1$ in total, all $\binom{N-1-k}{2}$ pairwise distances between the $N - k$ nodes in $\mathcal{L}$ are searched [38]. Using priority queue data structure we can reduce this complexity to $O(N^2 \log N)$. By using some more optimizations, it can be brought down to $O(N^2)$. In Matlab, hierarchical clustering implementation is usually $O(N^2)$. It is important to note

TABLE III
5G NR NUMEROLOGY AND SUBCARRIER SPACING [39]

| $\mu$ | Df= $2^\mu \cdot 15$ [kHz] | Bandwidth per RB [kHz] | TTI [ms] | Slots / ms |
|---|---|---|---|---|
| 0 | 15 | 180 | 1 | 1 |
| 1 | 30 | 360 | 0.5 | 2 |
| 2 | 60 | 720 | 0.25 | 4 |
| 3 | 120 | 1440 | 0.125 | 8 |
| 4 | 240 | 2880 | 0.0625 | 16 |

---

**Algorithm 2:** Optimal SA3M

1 **Input:** Multicast clusters $\mathcal{L}$;
2 Coordinates of $N$ multicast UEs $(X(i), Y(i), Z(i))$, $i \in \mathcal{N}$
3 **Output:** Optimal network configuration;
4 Create all $2^N$ possible network configurations considering unicasting;
5 Create all $2^N - 1$ possible network configurations considering D2D transmissions;
6 **for** *each network configuration* **do**
7      $T_{\text{total}}^{\text{NC}} = \sum_{m \in \mathcal{G}} T_m + \sum_{u \in \mathcal{U}} T_u$.
8 **end**
9 Solve optimization as per (12).

that these time complexities are general approximations and can vary based on factors such as the specific algorithm used, the distance metric employed, and the efficiency of the implementation.

### C. Step II – Optimization

Unicast, sidelink, and multicast transmission modes can coexist in a cell for the transmission of the same content. In this context, the UE tunes to the corresponding channel for data reception based on the optimization problem described below.

We consider a dynamic scenario where time is divided into discrete time slots $t$ of constant duration. 5G NR utilizes the scalable numerology that determines the subcarrier spacing, the number of slots in a subframe, and the slot duration (see Table III). At each time slot $t$, UEs can be associated with different transmission modes depending on the channel conditions, i.e., as per (5), (8), (10), and the results of the optimization.

The SA3M optimization deals with choosing the best network configuration in terms of the considered metric of interest (see Algorithm 2). To create the network configurations, the following rules are set: *(i)* UEs cannot join a multicast group different from the one defined using Algorithm 1 (we rerun MGF to form distinct multicast groups at a given rerunning interval, see Section VI-C); *(ii)* the predefined multicast transmission mode can be switched into unicast or sidelink for each UE following the SA3M algorithm to improve the network throughput.

As a preliminary step, the algorithm creates all possible network configurations that determine the transmission modes for all UEs (lines 4-5). Among them, the BS will choose (through exhaustive search) the one that optimizes the network performance. The number of possible network configurations is $2 \cdot 2^N - 1$ since there are $2^N$ possible combinations of 0 and 1, where 1 means that the UE remains in the multicast group determined by MGF, and 0 represents a switch of the

transmission mode to unicast. In the case of sidelink mode, we have $2^N - 1$ options as one of the devices should always be considered as a relay. Hence, all network configurations are in $2 \cdot 2^N - 1$.

Then, depending on the network configuration, UEs can be associated with different transmission modes, and SNR is determined as follows.

*Multicasting.* Multicast services are multi-user specific, and the quality of the channel of a multicast group $m$ is determined by the UE experiencing the worst channel conditions, i.e.,

$$S_m(t) = \min_{i \in \mathcal{G}_m} \left( \frac{P_T D_0 \rho(\alpha_i(t))}{N_0 W} \left[ \frac{y_i(t)^{-\zeta_{\mathrm{LoS}}}}{A_{\mathrm{LoS,nB}}(t)} \left[ 1 - p_B(y_i(t)) \right] \right. \right.$$
$$\left. \left. + \frac{y_i(t)^{-\zeta_{\mathrm{LoS}}}}{A_{\mathrm{LoS,B}}(t)} p_B(y_i(t)) \right] \right), \quad (5)$$

where $\mathcal{G}_m$ is the set of UEs in a multicast group covered by the same beam $m$, $1 < |\mathcal{G}_m| \le N$, $\mathcal{G}_m \subseteq \mathcal{N}$, $m \in \mathcal{G}$ is the subscript of the multicast group, and $\mathcal{G}$ is the set of all multicast groups. Initially, when all UEs are clustered into multicast groups, $\mathcal{L} = \mathcal{G}$, meaning that $\mathcal{L}$ is represented by multicast groups only. Set of unicast groups $\mathcal{U}$ is added at step 2.

The time required for the transmission of a packet of size $B$ to a multicast subgroup when experiencing the channel condition $S_m(t)$ can be calculated as

$$T_m(t) = \frac{B}{W_m \log_2(1 + S_m(t))}. \quad (6)$$

Hereinafter, we omit the slot notation $(t)$ for the sake of space.

The half power beamwidth (HPBW) $\theta$ required to serve subgroup $\mathcal{G}_m$ is given by:

$$\theta_{\mathcal{G}_m} = \arccos \left( \frac{X(i)X(i') + Y(i)Y(i') + Z(i)Z(i')}{y(i)y(i')} \right), \quad (7)$$

where multicast UEs $i$ and $i'$ are the two edge UEs in the group, i.e., the farthest in terms of the angle between them.

*Unicasting.* mmWave unicast transmission facilitates expanding the coverage area by sweeping narrow beams (e.g., HPBW of 2°). The UE that cannot be served as part of multicast transmission may prefer unicasting, considering the following link quality and data transmission duration:

$$S_u = \frac{P_T D_0}{N_0 W} \left( \frac{y_i^{-\zeta_{\mathrm{LoS}}}}{A_{\mathrm{LoS,nB}}} \left[ 1 - p_B(y_i) \right] + \frac{y_i^{-\zeta_{\mathrm{LoS}}}}{A_{\mathrm{LoS,B}}} p_B(y_i) \right), \quad (8)$$

$$T_u = \frac{B}{W_u \log_2(1 + S_u)}. \quad (9)$$

*D2D.* We assume in-band D2D, wherein UEs share the licensed uplink frequency resources with cellular communications. The channel link can be determined as

$$S_d = \frac{P_{T,d} D_0}{N_0 W} \left( \frac{y_{i,d}^{-\zeta_{\mathrm{LoS}}}}{A_{\mathrm{LoS,nB}}} \left[ 1 - p_B(y_{i,d}) \right] + \frac{y_{i,d}^{-\zeta_{\mathrm{LoS}}}}{A_{\mathrm{LoS,B}}} p_B(y_{i,d}) \right), \quad (10)$$

where $y_{i,d}$ is the distance between UE $i$ and D2D transmitter, and the data transmission delay can be calculated as

$$T_d = \frac{B}{W_d \log_2(1 + S_d)}. \quad (11)$$

The BS selects the possible D2D transmitter (relay device) based on the following rules: *(i)* the distance between a relay and a UE has to be within D2D$_{\mathrm{thr}}$ (i.e., $y_{i,d} <$ D2D$_{\mathrm{thr}}$), *(ii)* the closest relay among those that satisfy the previous condition is selected, and *(iii)* a relay can transmit data to one or more UEs (more details are given in the following). We assume that a relay device can simultaneously receive and transmit data to the UE (i.e., full-duplex relaying). A D2D transmission link cannot be established with a particular UE if no relay satisfies the described conditions. In this case, unicast transmission shall be performed.

We consider two relay selection options to account for different hardware on the devices. In the case referred to as "D2D communication without restrictions", a relay device can convey the traffic to more than one UE at a time. Conversely, a simple device works in the category of "D2D with restrictions", wherein it can relay data to one UE at a time.

We assume that the power transmitted by the relay node is lower than the power emitted by the BS, that is $P_{T,d} < P_T$, which helps to avoid that D2D communication causes exceeding interference. Recall that in-band D2D UEs reuse the same uplink resources of the mmWave cell, which can cause interference.

*1) Optimization Objective:* A multicast UE can receive data at different rates depending on its current location and blockage conditions. The optimization objective is to maximize the network throughput (NT) (i.e., aggregated throughput optimization). NT is calculated as the sum of data rates delivered to all UEs in the network.

Here, the problem consists of solving the overall maximum throughput optimization problem that is formulated as follows:

$$\max \quad \frac{BN}{\sum_{m \in \mathcal{G}} T_m + \sum_{u \in \mathcal{U}} T_u},$$
$$\text{s.t.} \quad S_m \ge S_{\mathrm{thr}}, S_u \ge S_{\mathrm{thr}}, S_d \ge S_{\mathrm{thr}}, y_{i,d} < y_{\mathrm{thr}}, \quad (12)$$

where $\mathcal{G}$ is the set of all multicast groups, and $\mathcal{U}$ is the set of all unicast users. As stated above, the algorithm needs to search through $2 \cdot 2^N - 1$ operations. Hence, the estimated complexity is $O(2^N)$. Note that $O(2^N)$ represents exponential time complexity, which means the computation time grows exponentially with the number of UEs, $N$. As $N$ increases, the computational requirements become increasingly prohibitive, and the exhaustive search becomes impractical. In such cases, heuristic algorithms or approximation techniques may be more suitable.

As an alternative to solving (12) according to SA3M, we propose a heuristic solution requiring low run-time, detailed in the following, to improve the system performance by adjusting the user-mode association.

*2) Proposed Heuristic:* The proposed heuristic, detailed in Algorithm 3, works as follows. First, it checks whether the SNR of each multicast group from $\mathcal{L}$ satisfies $S_{\mathrm{thr,h}}$ (line 6). If the SNR of the group (i.e., the worst SNR value among the group members as per (5)) is lower than the $S_{\mathrm{thr,h}}$ value, then the algorithm proceeds with checking every UE $i$ in this group (lines 7-8). In this case, UE $i$ is removed from the group and added to a separate one, $\mathcal{G}_g$, if its SNR value is below

**Algorithm 3:** Heuristic Solution

---

1   **Input:** Multicast clusters $\mathcal{L}$;
2   Coordinates of $N$ multicast UEs $(X(i), Y(i), Z(i)), i \in \mathcal{N}$
3   **Output:** Network configuration;
4   $g \leftarrow N$;
5   **for** *each $\mathcal{G}_m \in \mathcal{L}$* **do**
6     **if** $S_m < S_{thr,h}$ **then**
7       **for** *each UE $i \in \mathcal{G}_m$* **do**
8         **if** $S_i < S_{thr,h}$ **then**
9           $g \leftarrow g + 1$;
10           $\mathcal{L} \leftarrow \mathcal{L} \setminus \mathcal{G}_m$;
11           $\mathcal{G}_m \leftarrow \mathcal{G}_m \setminus \{x_i\}$;
12           $\mathcal{L} \leftarrow \mathcal{L} \cup \mathcal{G}_m$;
13           $\mathcal{G}_g \leftarrow \{x_i\}$;
14           $\mathcal{L} \leftarrow \mathcal{L} \cup \mathcal{G}_g$;

15   **for** *each $\mathcal{G}_m \in \mathcal{L}$* **do**
16     **if** $|\mathcal{G}_m| > 1$ **then**
17       calculate $\theta$ as per (7);
18       calculate $S_m, T_m$ as per (5),(6);
19     **else**
20       find $\max_{l \in \mathcal{G}^*}\{S(y_{i,l})|y_{i,l} < \mathrm{D2D}_{thr}\}$; $y_{i,l}$ is the distance between UEs $i$ and $l$, $\mathcal{G}^*$ is a set of UEs served via multicast
21       calculate $T_d$ as per (11);
22       calculate $S_u, T_u$ as per (8),(9);
23       choose best option as $\min(T_d, T_u)$;

24   **return** $\mathcal{L}$;
25   **end**

---

**TABLE IV**
DEFAULT PARAMETERS FOR NUMERICAL EVALUATION

| Parameter | Value |
|---|---|
| $f_c$ | 28 GHz |
| $W$ | 1 GHz |
| $R_d$ | 100 m [6], [28] |
| $h_A$ | 10 m |
| $h_U$ | 1.5 m |
| $h_B$ | 1.7 m |
| $r_B$ | 0.4 m |
| $\lambda_B$ | 0.3 bl/m² |
| $N_U$ | 4 el |
| $P_T$ | 46 dBm |
| $P_{T,d}$ | 10 dBm |
| $N_0$ | -174 dBm/Hz |
| $v$ | 0.69/11 m/sec |
| $S_{thr}$ | -9.47 dB |
| $S_{thr,h}$ | 6.367 dB (CQI 8) |
| $N$ | 10 [6], [28] |
| $M_{S,nB/B}$ | 4/8.2 dB |
| $B$ | 1 Gb |

the threshold (lines 6-14). By doing so, the algorithm detects the UEs that deteriorate the multicast group performance. The second *for* cycle of the algorithm is responsible for the calculation of the beamwidth $\theta$ of the multicast group (lines 16-19). All groups are already reformed at this stage, and the algorithm needs to adjust the beamwidth to be swept as per (7). Lines (19-23) are in charge of the selection between sidelink and unicast modes for single UEs that were deleted from multicast groups due to their channel conditions.

The computational complexity of the algorithm is given by

$$O((|\mathcal{L}| \cdot N) + |\mathcal{L}|) = O(|\mathcal{L}| \cdot N) = O\left(N^2\right),$$

where each summons (on the left side of the expression) determines the complexity of each *for* cycle. Then, as each cycle is called in turn (sequential execution), the complexity of the algorithm is $O(|\mathcal{L}| \cdot N)$. We note that $|\mathcal{L}| = N$ in the worst case when we have only unicast UEs. Hence, the algorithm's complexity is polynomial and can be rewritten as $O(N^2)$. Note that our algorithm has embedded *for* cycle (lines 7-14). We highlight that this execution helps reduce the complexity when not all multicast groups contain a "bad" UE that deteriorates the group's performance. These two cycles could be substituted by one *for* cycle among all $N$ UEs. That is, we could check all UEs without exclusion.

### D. Metrics of Interest

The analyzed metrics of interest are: *(i)* energy consumption, measured in joules (J), computed as the number of power units consumed over transmission time, i.e., in the case of multicasting, $EC = P_T \sum_{m \in \mathcal{G}} T_m$, *(ii)* network throughput representing the sum of data rates delivered to all UEs in the network as defined in (12), and *(iii)* latency, i.e.,

$\sum_{m \in \mathcal{G}} T_m$. We emphasize that the metrics are calculated for the resulting configuration after the execution of *(i)* both algorithms (MGF+SA3M) based on (12) and *(ii)* the proposed heuristic. We note that our ultimate goal is to maximize the network throughput. However, energy consumption is another critical metric to be considered in 5G/6G systems. Therefore, in the next section, we provide a set of results while carefully selecting the plotted metrics. That is, in case network throughput shows straightforward behavior, we present energy consumption plots instead and also analyze the impact of transmit power and available bandwidth on all metrics of interest.
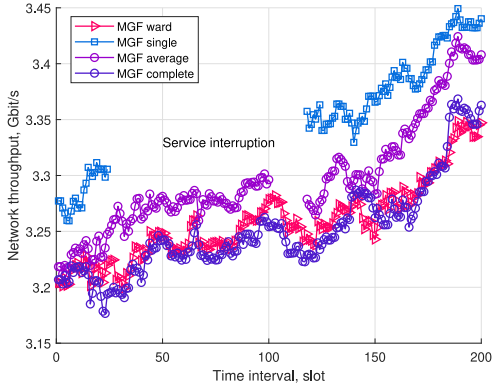
## VI. MAIN NUMERICAL RESULTS

This section describes the performance of the proposed MGF and SA3M algorithms and of the proposed heuristic solution, evaluated by means of a simulation environment developed in MATLAB that accepts the default parameters summarized in Table IV. In the remainder of the section, we first select the linkage function within the considered unsupervised learning algorithm that is best suited for MGF. We then proceed with a numerical analysis of the introduced optimal SA3M algorithm and discuss the effects of mobility, complexity, and UEs' distribution on the system performance and compare proposed algorithms against benchmarks. Finally, we report on the performance of the proposed low-complexity heuristic and analyze the impact of transmit power on energy consumption, latency, and network throughput.

Users are distributed within a sector of radius 100 m according to a Poisson point process (PPP) and Matérn cluster point process with 2 clusters. We consider $N = 10$ UEs, which is a commonly employed assumption in multicasting scenarios [28], [40]. The transmission parameters are modeled as indicated in Section IV with the operating frequency of 28 GHz and transmit power of 46 dBm. The bandwidth is 1 GHz [41] and the noise figure is 7.6 dB. The beam parameters are adjusted depending on the position of UEs in the multicast groups, whereas unicast and sidelink UEs utilize
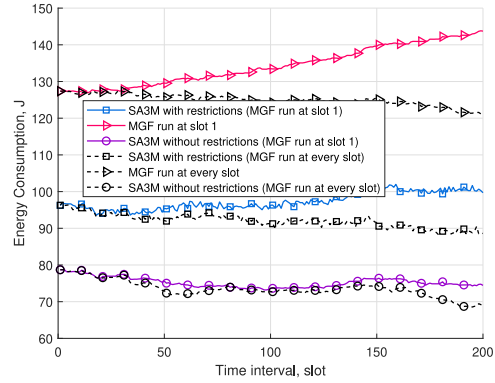
(a)



(b)

Fig. 3. Impact of linkage functions on (a) energy consumption and (b) network throughput over 200 time slots. Uniform user distribution.



(a) Uniform



(b) 2 clusters

Fig. 4. Energy consumption over time for pedestrian mobility: (a) uniform and (b) 2 clusters. Black lines are drawn in case MGF is rerun at every time slot.

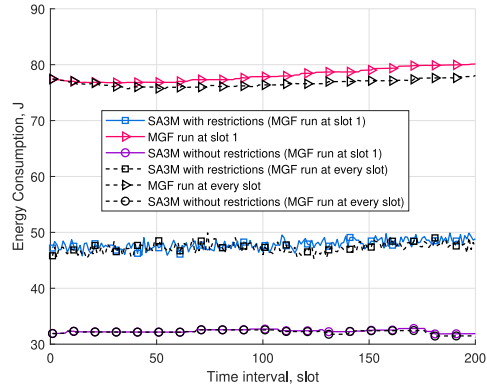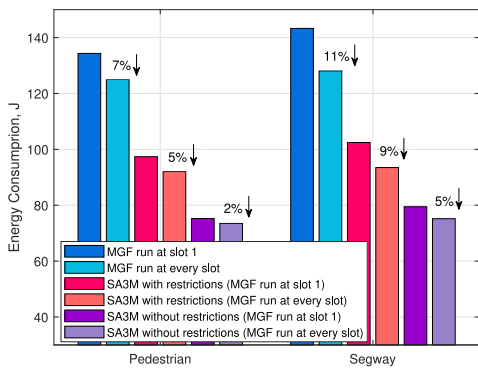the antenna with beamwidth of 3.18°. The mobility pattern of UEs is simulated as HSFM (see Section IV-D).

## A. Effect of Linkage Function

We test the four general linkage functions introduced in Section V-B for MGF to assess their impact on the performance. Our analysis focuses on their influence on energy consumption and network throughput, as depicted in Fig. 3. In Fig. 3(a), we observe a decreasing trend in the curves, indicating an improvement in energy consumption due to the MGF algorithm's ability to better track UEs mobility by rerunning it at every time slot. It is important to highlight that the *single* and *average* linkages exhibit better performance in terms of minimal energy consumption. However, they occasionally fail to maintain the minimum required service quality for all users of the multicast groups formed by the unsupervised learning algorithm. This can be attributed to group members being too far from each other. Consequently, the BS must utilize a wider beam to cover the entire group, resulting in a lower transmit antenna gain compared to more directional beam configurations. As a result, there might not be enough power to reach the worst user in the group.

Similar trends are observed in Fig. 3(b). In general, the *ward* function exhibits slightly superior performance in terms
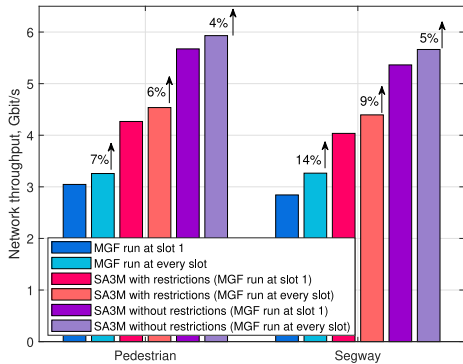
of mean energy consumption and network throughput compared to *complete* and significantly outperforms the other two functions (due to service interruption). This suggests that both ward and complete functions can capture directional multicast transmission features. Hence, we adopt the ward linkage as the default parameter for the MGF Algorithm in the following sections.

## B. Effect of Mobility

We start our primary evaluation campaign by analyzing the performance of the MGF and SA3M algorithms over time, as shown in Fig. 4. Our objective is to compare the performance of MGF and SA3M under different conditions, including restrictions for sidelink relaying, in two distinct modes: *(i)* MGF together with SA3M launched at every time slot and *(ii)* MGF launched at time slot 1 only (no rerunning) and SA3M run every time slot for the two considered user distributions with pedestrian mobility. We note that no rerunning of MGF (see colored curves compared to black ones) affects the performance over time, even though the speed and users' mobility are the same. In particular, for uniform distribution, we can see the most noticeable difference between the two running modes of the algorithms. In contrast, rerunning produces almost no improvement for the 2 cluster distribution, as seen in Fig. 4(b). This effect can be explained

(a) Energy consumption



(b) Network throughput

Fig. 5.  (a) Energy consumption and (b) network throughput for multicast users moving with different speeds in case of uniform distribution (each second bar is when we rerun MGF at every time slot).

by the fact that, in the case of PPP, UEs are spread around the area of interest. In general, the distance between every two uniformly distributed UEs in the network is higher than in the case of the cluster distribution. This impacts the performance since, in this case, D2D transmissions have to be performed over longer distances, and wider beams should be swept to cover multicast groups.

Analyzing further the effect of mobility for segway (with $v = 11\,\mathrm{m/sec}$) and pedestrian (with $v = 0.69\,\mathrm{m/sec}$) UEs in Fig. 5, we learn that for faster speeds, *MGF rerunning plays a crucial role in maintaining the performance level in dynamic scenarios*. As expected, the gap between rerunning vs. no rerunning of MGF is higher for segway mobility. The average performance improves by 11% and 7% for segway and walking UEs, respectively, in the case of energy consumption, and by 14% and 7% in the case of network throughput. Similarly, rerunning of MGF improves the performance of SA3M. Hence, it is highly recommended to rerun MGF to maintain the required performance level. In the following subsection, we comment on the rerunning interval of the MGF algorithm and on the complexity of the algorithms.

## C. Complexity vs. Energy Performance Trade-Off

We run the simulations via MATLAB R2021a on an Intel Core i5-7200U CPU @2.50 GHz at 2.71 GHz with 8.00 GB

TABLE V
ALGORITHMS' COMPLEXITY

| MGF | SA3M with restrictions | SA3M without restrictions | Heuristic |
|---|---|---|---|
| Execution time, seconds | | | |
| 0.0276 | 4.4580, of which: <br> - 3.8703 (line 4) <br> - 0.0022 (line 5) <br> - 0.5853 (lines 6-9) | 5.5605, of which: <br> - 3.8703 (line 4) <br> - 0.0022 (line 5) <br> - 1.6879 (lines 6-9) | 0.0313 |
| Theoretical complexity | | | |
| $O(N^2)$ | $O(2^N)$ | $O(2^N)$ | $O(N^2)$ |

RAM. The observed and theoretical complexities of the proposed algorithms are summarized in Table V.

We now analyze the complexity/energy performance trade-off in Fig. 6. For this reason, we run additional simulations to compare the performance of SA3M with or without MGF rerunning in terms of *(i)* energy consumption gain on the right *y*-axis and *(ii)* complexity gap on the left *y*-axis as a function of the MGF rerunning interval. First, let us analyze how MGF behaves when considering rerunning interval values ranging from 1 (i.e., the algorithm runs every single slot) to 40 (i.e., the algorithm runs every 40th slot). By observing Fig. 6(a), it emerges a high increase in complexity (up to 20000%) for MGF when rerunning is performed at every slot compared to the "no rerunning" case. Moreover, one may observe a noticeable drop in complexity for rerunning interval values ranging from 1 to 10. By further increasing the rerunning interval, the performance gap between "no rerunning" and "rerunning" slowly decreases.

Similar trends are observed in Fig. 6(b) and Fig. 6(c), in which, however, the observed quantitative increase in complexity is not so significant. This can be explained by the fact that the total complexity of SA3M is vastly greater than that of MGF. Hence, we may conclude that MGF complexity does not contribute to the overall complexity of SA3M. On the other hand, shortening the rerunning interval might be crucial for fast UE speeds, depending on the mobility pattern. Thus, taking into account both above-mentioned considerations, *our recommendation is to keep the rerunning interval in the range of* 10-20 *slots, which represents a good trade-off between complexity and achievable performance in terms of energy*.

## D. Effect of Users' Distribution and Benchmarks

This subsection illustrates the results of our evaluation campaign in terms of network throughput as a function of users' distribution. We compare the results of the proposed SA3M with several widely used methods: *(i)* sequential unicasting, *(ii)* D2D-assisted multicasting (based on Self-Organizing Map (SOM) unsupervised machine learning) [8], *(iii)* multicasting based on SOM with 4 neurons [8], *(iv)* incremental multicasting [30]. As shown in Fig. 7, the Matérn cluster distribution of UEs within the sector provides, in most cases, better results compared to the uniform one. The reason is that, in the case of Matérn cluster distribution, its randomly located points tend to be clusterized, which is beneficial for multicasting rather than being scattered around the area of interest. However, the observed trends in the results exhibit
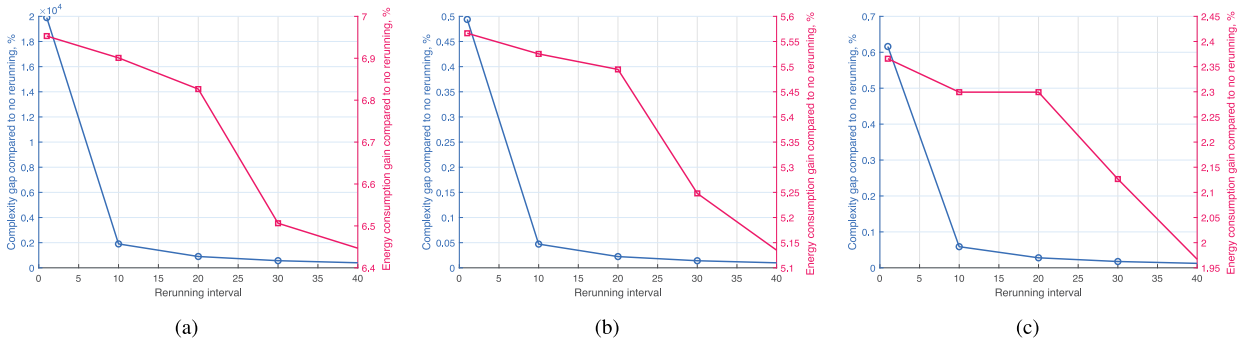
Fig. 6. Complexity gap vs. energy performance gain compared to no MGF rerunning for uniform distribution of UEs: (a) MGF, (b) SA3M without restrictions, (c) SA3M with restrictions.
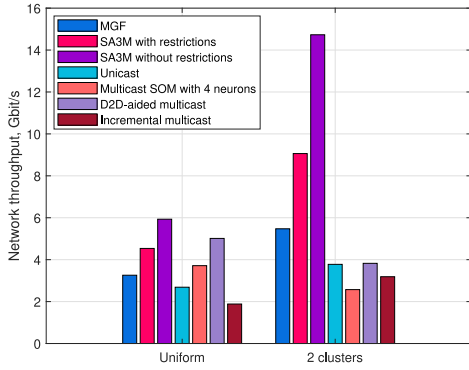


Fig. 7. Effect of different distributions of multicast users on network throughput in case MGF is executed at every time slot.
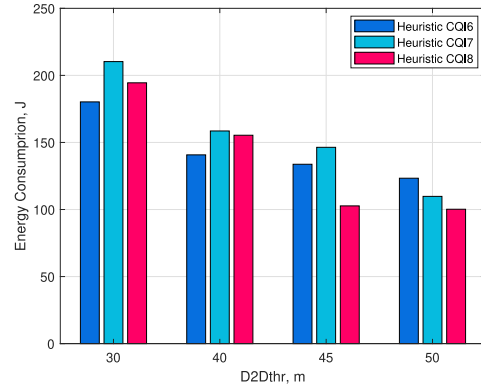


Fig. 8. Energy consumption for heuristic solution varying $S_{thr,h}$ and D2D$_{thr}$ thresholds.

differences depending on the distribution of multicast users, with certain benchmark schemes performing better for one distribution over another. For instance, "Multicast SOM with 4 neurons" outperforms "Incremental multicast" in the case of uniform distribution, while in the case of cluster distribution, the "Incremental multicast" emerges as the superior approach. Instead, the proposed SA3M algorithm without restrictions shows a dominant performance for both distribution types. This result highlights the versatility of the proposed solution and its potential for practical application.

### E. Heuristic Evaluation

To evaluate the performance of the heuristic algorithm, let us examine the impact of *(i) SNR threshold*, $S_{thr,h}$, for removing a multicast UE from the group (as for lines 6-14 of the Heuristic Algorithm) and assigning unicast/sidelink transmissions, and *(ii) distance threshold* at which sidelink communication can be established [42], D2D$_{thr}$. For SNR thresholds, we use modulation and coding scheme (MCS) mappings from [43].

The impact of D2D$_{thr}$ is evaluated in Fig. 8. The analysis demonstrates that using a less stringent threshold, specifically D2D$_{thr}$=50 m, for establishing sidelink communication distances outperforms all other thresholds. This superiority can be attributed to our proposed low-complexity heuristic, which avoids an exhaustive search for all possible configurations and instead generates configurations based on the channel conditions of UEs. In the case of thresholds such as 30 m, 40 m, and 45 m, uniformly distributed UEs are too far from each other and cannot fulfill the SNR requirements of the multicast group. Consequently, *bad* UEs must be removed from the multicast group and served via unicast links. This sequential transmission approach may degrade performance. Recall that we consider a single beam system and that higher CQI puts stricter requirements for multicast group channels. Hence, more sidelink connections should be established beyond multicast connections. This general trend, while varying the SNR threshold, is more remarkable for D2D$_{thr}$=50 m, where UEs can freely establish sidelink connections.

In Fig. 9, we further investigate the performance of the proposed Heuristic compared to MGF and both SA3M algorithms. In particular, we show the mismatch between the energy consumption required by the Heuristic with respect to the analyzed solutions. It can be noticed that the SNR threshold for CQI 8 provides the best Heuristic performance. The reason is that, in this case, more multicast users will be unclustered, and sidelink transmissions will be preferred. Furthermore, regarding the performance vs. complexity trade-off, it is essential to highlight that the Heuristic has lower complexity (by orders of magnitude) compared to SA3M while showing comparable performance under proper parameters' settings.
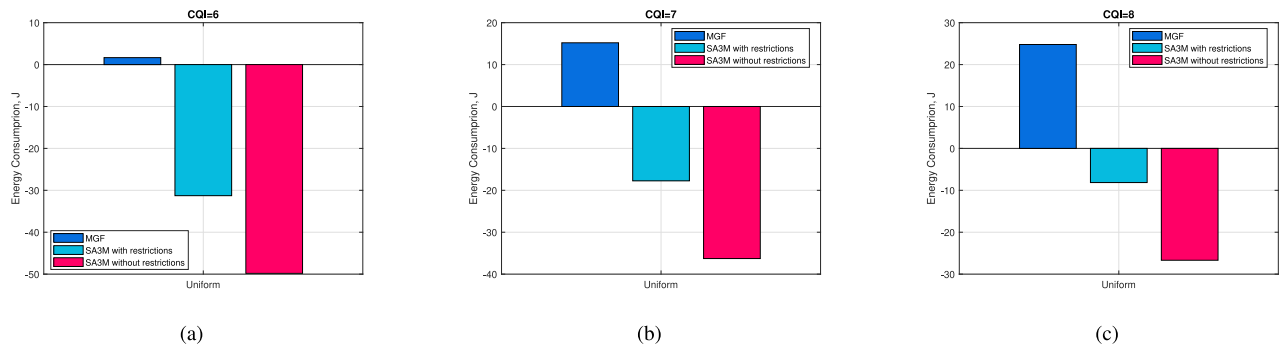
(a)                                                        (b)                                                        (c)

Fig. 9.   Heuristic compared to MGF and SA3M algorithms for (a) CQI 6, (b) CQI 7, and (c) CQI 8. $D2D_{thr}$=50 m. Effect of different distributions of multicast users on (a) energy consumption, (b) network throughput, and (c) energy efficiency. MGPF is executed at every time slot.



(a)                                                        (b)                                                        (c)



(d)                                                        (e)                                                        (f)
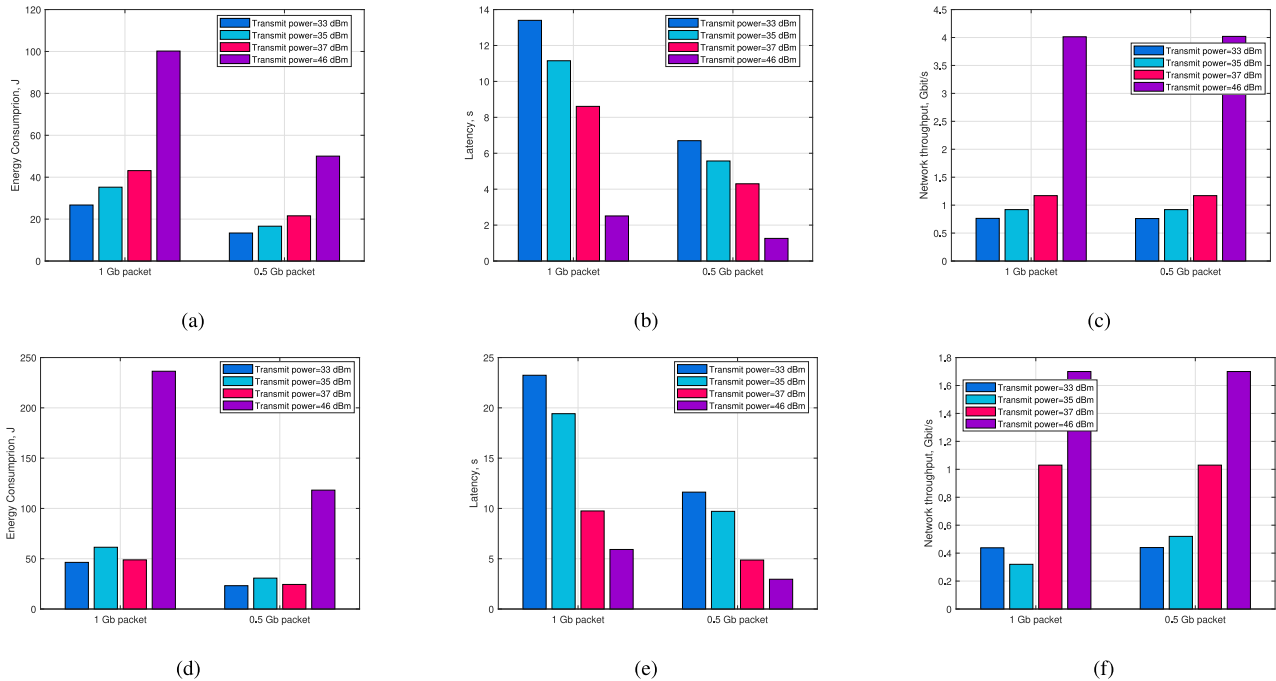
Fig. 10.   Effect of transmit power on (a and d) energy consumption, (b) latency, and (c) network throughput. Heuristic with $D2D_{thr}$=50 m, MCS 8, and $W = 1$ GHz (a, b, and c) and $W = 400$ MHz (d, e, and f).

## F. Effect of Transmit Power

As a final step, we investigate the impact of transmission power on energy consumption, latency, and network throughput. Obviously, by lowering the transmit power, a decrease in energy consumption is achieved. The same decreasing trend is experimented in terms of SNR, leading to a throughput degradation, which, in turn, causes a delay increase. The raised delay affects energy consumption. Therefore, there is a trade-off between transmit power, $P_T$, and delay. Moreover, packet size and available bandwidth at a transmission link also impact energy consumption.

To this aim, $P_T$ is varied as shown in Fig. 10(a,b,c). Observe that the rise in transmit power increases energy consumption up to 73% comparing the two extreme cases of 46 dBm and 33 dBm. Conversely, $P_T$ reduction results in higher latency, as demonstrated in Fig. 10(b). Depending on the service requirements and hardware on the devices, the choice of the transmit power can be shifted to one of the extreme cases or,

differently, to a middle value. By analyzing Fig. 10 further, one can notice that packet size does not influence the trend of the curves for $W = 1$ GHz since the available bandwidth in the system allows data delivery.

When focusing on Fig. 10(d,e,f), where the transmission bandwidth is set to $W = 400$ MHz, it is observed that the latency increases under decreasing transmission power values for all considered bandwidths and packet sizes. However, the performance in terms of energy consumption and network throughput shows different trends when varying the available bandwidth (see Fig. 10(a) w.r.t. Fig. 10(d), and Fig. 10(c) w.r.t. Fig. 10(f)). A nonlinear trend is observed in case of lower available bandwidth under increasing transmit power (see $P_T = 35$ dBm compared to $P_T = 33$ dBm and $P_T = 37$ dBm) for both considered packet sizes, revealing that there is a non-trivial relationship among transmit power, bandwidth, and latency and that the fine tuning of these parameters can lead to an advantageous reduction in both energy consumption and

latency. For example, see the magenta bars ($P_T = 37\,\text{dBm}$) in Fig. 10(d) w.r.t. the bars with other considered transmit powers and Fig. 10(a).

As a result, we infer that bandwidth, together with the transmit power, must be properly adjusted to reduce the total power consumption in the network.

## VII. CONCLUSION

In this work, we developed a two-step framework for sidelink-assisted multiple modes mmWave scheduling. Its complexity is reduced by exploiting an existing unsupervised learning algorithm to form multicast clusters. The resulting proposed solution leverages both optimization and machine learning techniques to deal with different types of users' mobility, user distribution, and network-side parameters, such as transmit power and bandwidth. To face complexity issues, a heuristic, which tracks channel conditions of multicast users, is also designed.

A thorough analysis of the system behavior has revealed crucial quantitative trade-offs to handle. Specifically, we elaborated on the complexity/performance trade-off connected with users' mobility. In particular, we recommend rerunning the multicast group formation algorithm, which allows for achieving better performance in mobile scenarios at the expense of low computational complexity. We then evaluated the energy consumption reduction as a consequence of a decreased transmit power and its impact on the total network latency. We emphasize that in 5G NR, the network's overall power consumption can be reduced by adjusting both bandwidth and transmit power. By combining the achieved results, we may conclude that multicast and D2D technologies are powerful tools to improve the performance of mmWave directional systems in the presence of dynamic users.

The findings of this research can be applied in the context of 5G networks, IoT deployments, smart city applications, vehicular communication systems, and other scenarios where reliable and high-performance wireless connectivity is crucial. By leveraging complexity reduction and adaptive resource allocation, these systems can enhance their efficiency, adaptability, and overall performance, thus meeting the demands of diverse communication environments.

Future research could explore the adaptation of the proposed framework for terahertz (THz) systems, where the coverage area of a single beam varies significantly depending on deployment specifics. This may involve incorporating advanced techniques such as ray-tracing to accurately model beam propagation and coverage patterns.

## REFERENCES

[1] M. Giordani, M. Polese, A. Roy, D. Castor, and M. Zorzi, "A tutorial on beam management for 3GPP NR at mmWave frequencies," *IEEE Commun. Surveys Tuts.*, vol. 21, no. 1, pp. 173–196, 1st Quart., 2019.

[2] N. Chukhno, O. Chukhno, S. Pizzi, A. Molinaro, A. Iera, and G. Araniti, "Approaching 6G use case requirements with multicasting," *IEEE Commun. Mag.*, vol. 61, no. 5, pp. 144–150, May 2023.

[3] N. Chukhno et al., "Models, methods, and solutions for multicasting in 5G/6G mmWave and sub-THz systems," *IEEE Commun. Surveys Tuts.*, vol. 26, no. 1, pp. 119–159, 1st Quart., 2024.

[4] S. Pizzi, C. Suraci, A. Iera, A. Molinaro, and G. Araniti, "A sidelink-aided approach for secure multicast service delivery: From human-oriented multimedia traffic to machine type communications," *IEEE Trans. Broadcast.*, vol. 67, no. 1, pp. 313–323, Mar. 2021.

[5] N. Chukhno, A. Orsino, J. Torsner, A. Iera, and G. Araniti, "5G NR sidelink multi-hop transmission in public safety and factory automation scenarios," *IEEE Netw.*, vol. 37, no. 5, pp. 129–136, Sep. 2023.

[6] A. Biason and M. Zorzi, "Multicast transmissions in directional mmWave communications," in *Proc. 23th Eur. Wireless Conf.*, 2017, pp. 1–7.

[7] N. Chukhno, O. Chukhno, S. Pizzi, A. Molinaro, A. Iera, and G. Araniti, "Efficient management of multicast traffic in directional mmWave networks," *IEEE Trans. Broadcast.*, vol. 67, no. 3, pp. 593–605, Sep. 2021.

[8] N. Chukhno, O. Chukhno, S. Pizzi, A. Molinaro, A. Iera, and G. Araniti, "Unsupervised learning for D2D-assisted multicast scheduling in mmWave networks," in *Proc. IEEE BMSB*, 2021, pp. 1–6.

[9] Y. Niu, Y. Liu, Y. Li, X. Chen, Z. Zhong, and Z. Han, "Device-to-device communications enabled energy efficient multicast scheduling in mmWave small cells," *IEEE Trans. Commun.*, vol. 66, no. 3, pp. 1093–1109, Mar. 2018.

[10] Y. Niu, L. Yu, Y. Li, Z. Zhong, and B. Ai, "Device-to-device communications enabled multicast scheduling for mmWave small cells using multi-level codebooks," *IEEE Trans. Veh. Technol.*, vol. 68, no. 3, pp. 2724–2738, Mar. 2019.

[11] Y. Niu, L. Yu, Y. Li, Z. Zhong, B. Ai, and S. Chen, "Device-to-device communications enabled multicast scheduling with the multi-level codebook in mmWave small cells," *Mobile Netw. Appl.*, vol. 24, no. 5, pp. 1603–1617, 2019.

[12] G. H. Sim, M. Mousavi, L. Wang, A. Klein, and M. Hollick, "Joint relaying and spatial sharing multicast scheduling for mmWave networks," in *Proc. IEEE 21st Int. Symp. World Wireless, Mobile Multimedia Netw. (WoWMoM)*, 2020, pp. 127–136.

[13] S. Zhang, D. Liu, J. Lv, and Z. Zhang, "D2D-enabled multicast optimal scheduling in mmWave cellular networks," in *Proc. IEEE/CIC Int. Conf. Comm. China (ICCC)*, 2020, pp. 442–447.

[14] J. Dai, G. Yue, S. Mao, and D. Liu, "Sidelink-aided multiquality tiled 360° virtual reality video multicast," *IEEE Internet Things J.*, vol. 9, no. 6, pp. 4584–4597, Mar. 2022.

[15] T. Soni, M. Schellmann, J. Eichinger, and A. Knoll, "Unleashing latency-critical IIoT communication by virtue of cooperative sidelink-assisted DL transmissions," in *Proc. 25th Int. ITG Workshop Smart Antennas (WSA)*, 2021, pp. 1–6.

[16] S. Naribole and E. Knightly, "Scalable Multicast in highly-directional 60-GHz WLANs," *IEEE/ACM Trans. Netw.*, vol. 25, no. 5, pp. 2844–2857, Oct. 2017.

[17] H. D. R. Albonda and J. Pérez-Romero, "An efficient mode selection for improving resource Utilization in sidelink V2X cellular networks," in *Proc. IEEE Int. Workshop Comput. Aided Model. Design Commun. Links Netw. (CAMAD)*, 2018, pp. 1–6.

[18] R. W. Heath, N. González-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, "An overview of signal processing techniques for Millimeter wave MIMO systems," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 3, pp. 436–453, Apr. 2016.

[19] L. Feng, Z. Yang, Y. Yang, X. Que, and K. Zhang, "Smart mode selection using online reinforcement learning for VR broadband broadcasting in D2D assisted 5G HetNets," *IEEE Trans. Broadcast.*, vol. 66, no. 2, pp. 600–611, Jun. 2020.

[20] S. Aslam, F. Alam, S. F. Hasan, and M. A. Rashid, "A machine learning approach to enhance the performance of D2D-enabled clustered networks," *IEEE Access*, vol. 9, pp. 16114–16132, 2021.

[21] Y. Wang, M. Narasimha, and R. W. Heath, "MmWave beam prediction with situational awareness: A machine learning approach," in *Proc. IEEE 19th SPAWC*, 2018, pp. 1–5.

[22] G. H. Sim, S. Klos, A. Asadi, A. Klein, and M. Hollick, "An online context-aware machine learning algorithm for 5G mmWave vehicular communications," *IEEE/ACM Trans. Netw.*, vol. 26, no. 6, pp. 2487–2500, Dec. 2018.

[23] F. GÖttsch and M. Kaneko, "Deep learning-based beamforming and blockage prediction for sub-6GHz/mmWave mobile networks," in *Proc. IEEE Globecom*, 2020, pp. 1–6.

[24] X. Wang et al., "Mean field reinforcement learning based anti-jamming communications for ultra-dense Internet of Things in 6G," in *Proc. Int. Conf. Wireless Comm. Signal Process. (WCSP)*, 2020, pp. 195–200.

[25] H. Wei, G. Zheng, V. Gayah, and P. Li, "Recent advances in reinforcement learning for traffic signal control: A survey of models and evaluation," *ACM SIGKDD Explor. Newslett.*, vol. 22, no. 2, pp. 12–18, 2021.

[26] D. Lee, N. He, P. Kamalaruban, and V. Cevher, "Optimization for reinforcement learning: From a single agent to cooperative agents," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 123–135, May 2020.

[27] H. Liu and W. Wu, "Federated reinforcement learning for Decentralized voltage control in distribution networks," *IEEE Trans. Smart Grid*, vol. 13, no. 5, pp. 3840–3843, Sep. 2022.

[28] N. Chukhno et al., "Optimal multicasting in Millimeter wave 5G NR with multi-beam directional antennas," *IEEE Trans. Mobile Comput.*, vol. 22, no. 6, pp. 3572–3588, Jun. 2023.

[29] H. H. Park and C.-H. Kang, "A group-aware multicast scheme in 60GHz WLANs," *Trans. Internet Inf. Syst.*, vol. 5, no. 5, pp. 1028–1048, 2011.

[30] H. Park, S. Park, T. Song, and S. Pack, "An incremental multicast grouping scheme for mmWave networks with directional antennas," *IEEE Commun. Lett.*, vol. 17, no. 3, pp. 616–619, Mar. 2013.

[31] S. Pizzi, M. Condoluci, G. Araniti, A. Molinaro, A. Iera, and G.-M. Muntean, "A unified approach for efficient delivery of unicast and multicast wireless video services," *IEEE Trans. Wireless Comm.*, vol. 15, no. 12, pp. 8063–8076, Dec. 2016.

[32] A. Samuylov et al., "Characterizing resource allocation trade-offs in 5G NR serving multicast and unicast traffic," *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3421–3434, May 2020.

[33] O. Chukhno et al., "A holistic assessment of directional deafness in mmWave-based distributed 3D networks," *IEEE Trans. Wireless Commun.*, vol. 21, no. 9, pp. 7491–7505, Sep. 2022.

[34] "Study on channel model for frequencies from 0.5 to 100 GHz; (Release 14), Version 14.1.1," 3GPP, Sophia Antipolis, France, Rep. TR 38.901, Jul. 2017.

[35] M. Gapeyenko et al., "Analysis of human-body blockage in urban Millimeter-wave cellular communications," in *Proc. IEEE Int. Conf. Commun. (ICC)*, 2016, pp. 1–7.

[36] F. Farina, D. Fontanelli, A. Garulli, A. Giannitrapani, and D. Prattichizzo, "Walking ahead: The headed social force model," *PloS One*, vol. 12, no. 1, 2017, Art. no. e0169734.

[37] F. Nielsen, "Hierarchical clustering," in *Introduction to HPC with MPI for Data Science*. Berlin, Germany: Springer, 2016, pp. 195–211.

[38] D. Müllner, "Modern hierarchical, agglomerative clustering algorithms," 2011, *arXiv:1109.2378*.

[39] "NR; physical channels and modulation; (Release 15)," 3GPP, Sophia Antipolis, France, Rep. TR 38.211, Dec. 2017.

[40] A. Biason and M. Zorzi, "Multicast via point to multipoint transmissions in directional 5G mmWave communications," *IEEE Commun. Mag.*, vol. 57, no. 2, pp. 88–94, Feb. 2019.

[41] S. Akoum, O. El Ayach, and R. W. Heath, "Coverage and capacity in mmWave cellular systems," in *Proc. 46th Conf. Rec. Asilomar Conf. Signals, Syst. Comput. (ASILOMAR)*, 2012, pp. 688–692.

[42] R. Li, P. Hong, K. Xue, M. Zhang, and T. Yang, "Energy-efficient resource allocation for high-rate underlay D2D communications with statistical CSI: A one-to-many strategy," *IEEE Trans. Veh. Technol.*, vol. 69, no. 4, pp. 4006–4018, Apr. 2020.

[43] J. Fan, Q. Yin, G. Y. Li, B. Peng, and X. Zhu, "MCS selection for throughput improvement in downlink LTE systems," in *Proc. 20th Int. Conf. Comput. Commun. Netw. (ICCCN)*, 2011, pp. 1–5.

**Olga Chukhno** (Member, IEEE) received the B.Sc. degree in business informatics from RUDN University, Russia, in 2017, the M.Sc. degree in fundamental informatics and information technologies in 2019, and the Double Ph.D. degrees in information engineering with the H2020 MCSA ITN/EJD Project, Mediterranea University of Reggio Calabria, Italy, and Tampere University, Finland. She is an Assistant Professor in telecommunications with the University Mediterranea of Reggio Calabria and CNIT, Italy. Her research interests include wireless communications and edge computing.

**Sara Pizzi** (Member, IEEE) received the first- and second- level Laurea degree (cum laude) in telecommunication engineering, the master's degree in information technology from CEFRIEL/Politecnico di Milano in 2005, and the Ph.D. degree in computer, biomedical and telecommunication engineering from the University Mediterranea of Reggio Calabria, Italy, 2002 and 2005, and 2009, respectively, where she is an Assistant Professor of Telecommunications, with the University Mediterranea of Reggio Calabria, Italy and CNIT, Italy. Her research interests focus on resource management, multicasting, D2D and MTCs over 5G/6G networks, and NTN.

**Antonella Molinaro** (Senior Member, IEEE) received Graduation in computer engineering from the University of Calabria in 1991, the master's degree in information technology from CEFRIEL/Polytechnic of Milano in 1992, and the Ph.D. degree in multimedia technologies and communications systems in 1996. She is currently a Full Professor of Telecommunications with the University Mediterranea of Reggio Calabria, Italy, CNIT, Italy, and Université Paris–Saclay, France. Her research activity mainly focuses on wireless and mobile networking, vehicular networks, and future Internet.

**Antonio Iera** (Senior Member, IEEE) received Graduation in computer engineering from the University of Calabria in 1991, and the master's degree in information technology from the CEFRIEL/Politecnico di Milano in 1992, and the Ph.D. degree from the University of Calabria in 1996. From 1997 to 2019, he has been with the University Mediterranea, Italy, and currently holds the position of Full Professor of telecommunications with the University of Calabria, Italy and CNIT, Italy. His research interests include next generation mobile and wireless systems, and the Internet of Things.

**Nadezhda Chukhno** (Member, IEEE) received Graduation from RUDN University, Russia, the B.Sc. degree in business informatics in 2017, the M.Sc. degree in fundamental informatics and information technologies in 2019, and the Double Ph.D. degree in information engineering with the H2020 MCSA ITN/EJD A-WEAR Project, Mediterranea University of Reggio Calabria, Italy and Jaume I University, Spain, in 2023. She is a Postdoctoral Reseacher with Tampere University, Finland. Her research interests include wireless communications, 5G+ networks, and ML.

**Giuseppe Araniti** (Senior Member, IEEE) received the Laurea degree and the Ph.D. degree in electronic engineering from the University Mediterranea of Reggio Calabria, Italy, in 2000 and 2004, respectively, where he is currently an Associate Professor of Telecommunications with the University Mediterranea of Reggio Calabria, Italy, and CNIT, Italy. His major area of research is on 5G/6G networks and it includes personal communications, wireless and satellite systems, traffic and radio resource management, multicast and broadcast services, and D2D.