# DM2-ECOP: An Efficient Computation Offloading Policy for Multi-user Multi-cloudlet Mobile Edge Computing Environment

HOUSSEMEDDINE MAZOUZI and NADJIB ACHIR, L2TI, Institut Galilée, Université Paris 13, Sorbonne Paris Cité, France

KHALED BOUSSETTA, L2TI, Institut Galilée, Université Paris 13, Sorbonne Paris Cité, Agora/INRIA, France

Mobile Edge Computing is a promising paradigm that can provide cloud computing capabilities at the edge of the network to support low latency mobile services. The fundamental concept relies on bringing cloud computation closer to users by deploying cloudlets or edge servers, which are small clusters of servers that are mainly located on existing wireless Access Points (APs), set-top boxes, or Base Stations (BSs). In this article, we focus on computation offloading over a heterogeneous cloudlet environment. We consider several users with different energy—and latency-constrained tasks that can be offloaded over cloudlets with differentiated system and network resources capacities. We investigate offloading policies that decide which tasks should be offloaded and select the assigned cloudlet, accordingly with network and system resources. The objective is to minimize an offloading cost function, which we defined as a combination of tasks' execution time and mobiles' energy consumption. We formulate this problem as a Mixed-Binary Programming. Since the centralized optimal solution is NP-hard, we propose a distributed linear relaxation-based heuristic approach that relies on the Lagrangian decomposition method. To solve the subproblems, we also propose a greedy heuristic algorithm that computes the best cloudlet selection and bandwidth allocation following tasks' offloading costs. Numerical results show that our offloading policy achieves a good solution quickly. We also discuss the performances of our approach for large-scale scenarios and compare it to state-of-the-art approaches from the literature.

CCS Concepts: • **Networks → Cloud computing**; • **Human-centered computing → Ubiquitous and mobile computing systems and tools**; • **Theory of computation → Network optimization**;

Additional Key Words and Phrases: Computation offloading, mobile cloud computing, mobile edge computing, cloudlet, Lagrangian decomposition

Authors' addresses: H. Mazouzi and N. Achir, L2TI, Institut Galilée, Université Paris 13, Sorbonne Paris Cité, 99 Avenue J-B Clement, Villetaneuse, 93430, France; emails: {mazouzi.houssemeddine, nadjib.achir}@univ-paris13.fr; K. Boussetta, L2TI, Institut Galilée, Université Paris 13, Sorbonne Paris Cité, Agora/INRIA, 99 Avenue J-B Clement, Villetaneuse, 93430, France; email: khaled.boussetta@univ-paris13.fr.

ACM Transactions on Internet Technology, Vol. 19, No. 2, Article 24. Publication date: March 2019.

**24**

## 1    INTRODUCTION

The Mobile Cloud Computing paradigm has been proposed to allow remote execution of resource-hungry mobile applications in the cloud. The application's computation is then transmitted to the remote cloud to be performed. The latter operation is known as **computation offloading** [6, 7]. Unfortunately, the geographical distance between the cloud and user can introduce large and variable latency. That can significantly degrade the quality of experience of delay-sensitive applications, such as mobile gaming, augmented-reality, and face and speech recognition [8, 30]. To overcome such problems, **Mobile Edge Computing (MEC)** has emerged as a main paradigm that aims to provide cloud computing capabilities at the edge of the network to support latency-sensitive mobile applications. The main concept relies on deploying small clusters of servers, called cloudlets, at the edge of the network [11, 18]. Users can then offload their computation to closer cloudlets.

In multi-user context, several mobile devices can compete to offload their computations to the cloudlets. Hence, the performances of offloading policies are strongly dependent on the cloudlets' computational resources sharing and on the wireless bandwidth allocation strategies [4, 5, 12]. In addition, in a multi-cloudlet MEC environment, where many cloudlets are available around users, the performance of the computation offloading depends on the cloudlet selection [14, 33, 34].

Many recent works have investigated cloudlet selection problems [21, 27, 34]. Most of the proposed offloading policies rely on user density to statically assign each region to a cloudlet [15, 27]. Therefore, as shown in Figure 1, users within a region will always offload to the same cloudlet. Nevertheless, the dynamic density of users may imbalance the load between the cloudlets, leading to suboptimal MEC capacities usage and longer offloading delays. Therefore, to achieve high performance, an offloading policy must jointly consider bandwidth allocation, computation resource allocation, and cloudlet selection.

To tackle this problem, we explore, in this article, computation offloading in multi-user, multi-cloudlet MEC. Our aim is to provide an efficient offloading policy, which determines the best offloading decision and cloudlet selection for each user, with the aim of reducing the total offloading cost.

This work presents a new computation offloading policy named Distributed Multi-user Multi-cloudlet Efficient Computation Offloading Policy (DM2-ECOP), which aims to improve the performance of offloading in an MEC environment. It extends the offloading strategies presented in References [4, 5, 12]. As in previous works, we assume that each user executes only one application at a time. However, unlike these works, we define two categories of applications that can be supported by the MEC: (1) applications that must be performed remotely and (2) applications that can be performed either locally or remotely, accordingly with the conditions at execution time. DM2-ECOP tries to select the best cloudlet according to network and system resource availability, while minimizing the offloading cost. The offloading cost is defined as a combination of the energy consumed by the mobile devices and the total applications' completion times.

We formulate this computation offloading problem as a Mixed Binary Programming. Then, we solve it using a distributed linear relaxation-based heuristic that follows the Lagrangian decomposition approach. DM2-ECOP is composed of two decision levels: (1) The local offloading manager handles the users associated within the same AP and solves the offloading subproblem related to this AP; the local offloading manager uses our proposed algorithm, named Greedy Best Cloudlet Selection First Heuristic (GBC-SFH), which selects the cloudlet to which each application will be offloaded to minimize the energy consumption and completion times. (2) At a second level, a global offloading manager ensures that the cloudlets' resources allocated by each local offloading manager satisfy the capacity constraints of each cloudlet.
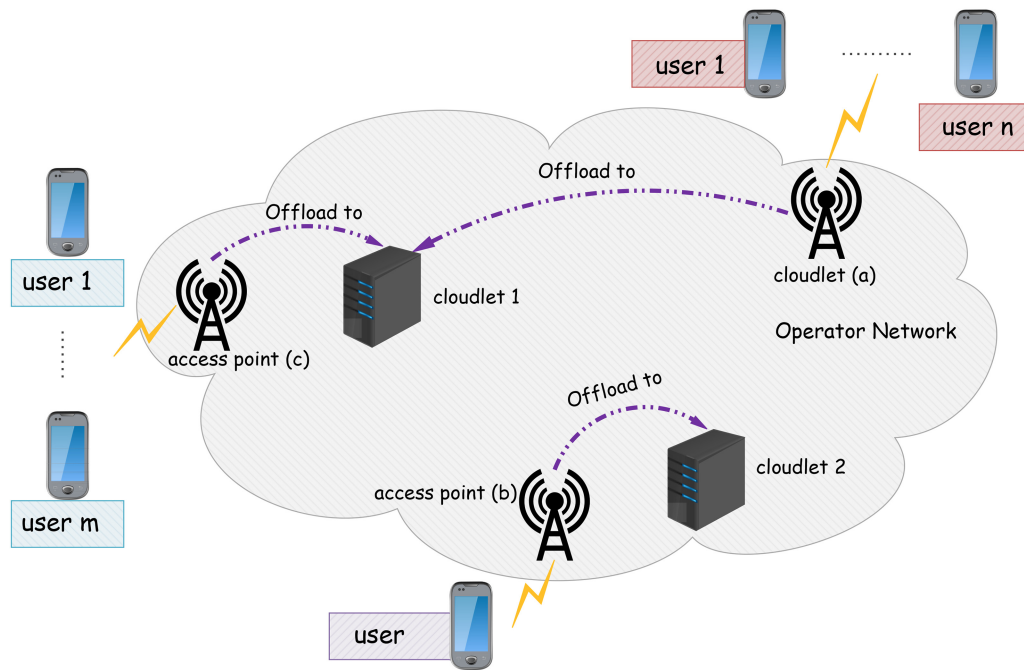
Fig. 1.  Multiuser computation offloading in multi-cloudlet MEC environment.

The remainder of this article is organized as follows: Section 2 presents existing works in computation offloading. Section 3 introduces the system modeling. The multi-user, multi-cloudlet offloading problem is formulated in Section 4. Our offloading policy, named DM2-ECOP, is explained in Section 5. Performance evaluation is detailed and analyzed in Section 6. Finally, a conclusion is drawn in Section 7.

## 2   RELATED WORK

Many works were proposed to explore computation offloading to improve the performance of mobile devices. However, not all of the proposed offloading policies have the same goals. In the following, we distinguish between three main goals: (i) offloading decision, (ii) cloudlets placement, and (iii) cloudlet selection.

Some of the proposed works investigate the offloading decision to decide which computation should be offloaded to the remote cloud, such as: Meng-Hsi Chen et al. [4], Xu Chen et al. [5], Songtao Guo et al. [12], Keke Gai et al. [10], Yuyi Mao et al. [22], and Dong Huang et al. [13]. Meng-Hsi et al. are one of the first to work on multi-user computation offloading in mobile cloud computing. The proposed offloading policy determines which computation must be performed in the remote cloud and which one must be performed locally by the mobile device. Then, it allocates the wireless bandwidth to each user to reduce the energy consumption of the mobile device. The Xu Chen et al. offloading policy was designed for a single cloudlet MEC environment. Each user tries to offload its computation, accordingly with the available wireless bandwidth to reduce the energy consumption. Another offloading approach for multi-user was presented by Songtao Guo et al. Similar to DM2-ECOP, this work minimizes an offloading cost defined as a combination of energy consumption and processing time. The offloading policy decides which computation can be offloaded and allocates the wireless bandwidth and the processor frequency to each offloaded

100 computation. In similar way, Keke Gai et al. proposed a scheduler to assign the tasks between the
101 local mobile device and the remote cloud to save energy consumption. Yuyi Mao et al. presented an
102 offloading policy that tries to offload the computation in a multi-user scenario to an edge server,
103 cloudlet. The proposed policy allocates the CPU frequency and the bandwidth to each user to
104 reduce the energy consumption of the mobile device. Lastly, Dong Huang et al. designed a compu-
105 tation offloading policy for a single-user scenario to reduce the energy consumption of the mobile
106 device. They focused on partial offloading, where the offloading policy partitions the application at
107 runtime to determine which computation must be performed locally and which must be offloaded
108 to a remote server. Although all these policies improve the performances of the mobile device,
109 they rely on an unlimited capacity of the cloud. Consequently, they need some enhancements to
110 be applied for the MEC, where cloudlets have limited computing resources.

111 Yucen Nan et al. [25, 26] and Chongyu Zhou et al. [35] proposed computation offloading poli-
112 cies to reduce the energy consumption of fog nodes. They introduced an offloading policy where
113 the fog nodes try to offload their computation to the remote cloud. For each fog node, the policy
114 decides which computation must be offloaded to the remote cloud and which one must be per-
115 formed locally by the node. In Reference [26], the offloading policy has been extended to reduce
116 the completion time of the applications. Similarly, Chongyu Zhou et al. introduced an online of-
117 floading policy. It can select the computations that should be performed by the nearest cloudlet
118 in order to minimize a system-wide utility, which is the execution time. Contrary to these poli-
119 cies, which reduce the energy consumption of the fog server, DM2-ECOP focuses on the reduction
120 of the offloading cost on the mobile device side. In addition, the IoT device has a tiny comput-
121 ing capacity that cannot perform any application. However, the mobile device has considerable
122 computing capacity that performs complex applications.

123 Cloudlets placement is also a challenging issue for MEC, and many recent works propose some
124 cloudlets placement heuristics in MEC environment. Mike Jia et al. [14, 33], Hong Yao et al. [34],
125 and Longjie Ma et al. [21] introduced cloudlets placement and selection algorithms in a multi-user,
126 multi-cloudlet MEC environment. The Mike Jia et al. offloading approach is one of the first heuris-
127 tics on cloudlets placement in a large-scale environment. Its main goal is to find the best cloudlets
128 placement in a large network, then select a cloudlet to perform the computation of each AP. The
129 K-median clustering based on user density is used to place the cloudlets. Then each AP is stati-
130 cally assigned to a cloudlet. Similarly, Hong Yao et al. have been designing heuristics to support
131 heterogeneous cloudlets environment. Finally, Longjie Ma et al. have been introducing a heuris-
132 tic to find the minimal number of cloudlets that must be placed to improve the user experience
133 quality in a large-scale network. In a multi-user MEC environment, the density of mobile users is
134 dynamic and changes over time. So, static assignment of the APs to cloudlets may decrease the
135 performance of the computation offloading. To avoid this problem, our DM2-ECOP approach con-
136 siders dynamic cloudlet selection and wireless bandwidth allocation with the aim of minimizing
137 energy consumption and improving the performance of mobile devices.

138 Other works try to find a dynamic cloudlet selection in a multi-cloudlet MEC environment.
139 Anwesha Mukherjee et al. [24, 27], Mike Jia et al. [15], Qiliang Zhu et al. [36], and Arash
140 Bozorgchenani et al. [2] have proposed to support the dynamic cloudlet selection to reduce the of-
141 floading cost. Anwesha Mukherjee et al. designed a multilevel offloading policy to optimize energy
142 consumption. The users offload to the nearest cloudlet in the first step. According to the amount of
143 resources available in this cloudlet, it can perform the task or offload it to another cloudlet. Mike
144 Jia et al. introduced a heuristic to balance the load between the cloudlet. Its main goal is to migrate
145 some computations from overloaded cloudlets to underloaded cloudlets to reduce the execution
146 time. Similarly, Qiliang Zhu et al. developed a two-tier offloading policy, where the mobile device
147 offloads its computation to an offloading server based on the resource availability. They used an

agent that decides to perform the computation in the local cloudlet or to offload it to the remote    148
cloud. Arash Bozorgchenani et al. offloading policy tries to select a nearby fog node to offload    149
some computation of a busy fog node to save energy consumption and completion time. Even    150
these works proposed dynamic cloudlet selection heuristics, the tasks still always offloaded to the    151
nearest cloudlet that decides to perform them locally or transmit them to other cloudlets. Thus, an    152
additional offloading cost is induced; consequently, the performance of offloading will decrease. In    153
our proposal, the cloudlet that performs each offloaded task will be determined at the offloading    154
decision time without any additional cost.    155

All the offloading policies presented above have been focusing on reducing the offloading cost.    156
They offload computations to a predetermined remote server (the remote cloud or a local cloudlet).    157
The selection of the remote server is done statically at the development time based on metrics such    158
as user density, despite the fact that the density of users can change dynamically over the time.    159
In addition, the computing capacity of the cloudlet is limited and cannot perform all the offloaded    160
computation. To avoid this situation, the most adopted strategy in the literature was to consider    161
a two-tier approach. Basically, the tasks are offloaded to the nearest cloudlet, and this cloudlet    162
offloads some computation to another cloudlet or to the remote cloud when it is overloaded.    163
Although two-tier offloading policies can improve the performance of the offloading approach,    164
they engender an additional offloading cost. Moreover, in a multi-cloudlet scenario, where many    165
cloudlets are available around the user, selecting the same cloudlet always is not the best strategy.    166
Therefore, in this article, we propose a new offloading policy to improve the efficiency of com-    167
putation offloading in MEC. The new policy must consider many cloudlets for which a user can    168
offload its computation and compute optimal computation placements to optimize the offloading    169
cost.    170

## 3  SYSTEM DESCRIPTION    171

In this section, we describe our system modeling. We first introduce the MEC model, then we    172
present the communication and computation offloading models. Finally, an offloading cost is pro-    173
posed as an objective function for our optimization problem. Table 1 presents variables and nota-    174
tions used, in this article, to model our multi-user, multi-cloudlet computation offloading problem.    175

### 3.1  MEC Environment Model    176

Let us consider an MEC environment composed of $M$ APs and $K$ cloudlets, as illustrated in Figure 2.    177
We suppose that the number of cloudlets is less than the number of APs ($K \leq M$). In this article,    178
we assume that the cloudlets have already been deployed and are co-located with the APs. We    179
also consider that the users can communicate with the cloudlets through their APs. We denote    180
in the following the set of APs by $\mathcal{M} = \{1, 2, \ldots, M\}$, and we assume that each AP $i$ is associated    181
with $N_i$ users. Let us consider $\mathcal{N}_m = \{1, 2, \ldots, N_m\}$ as the set of users associated with the $m^{th}$ AP    182
and $\mathcal{K} = \{1, 2, \ldots, K\}$ as the set of the cloudlets. We also define $u_{m,n}$ as to $n^{th}$ user of the $m^{th}$ AP.    183
Similar to existing works [14, 22, 33], every user runs one application on his mobile device. The    184
application is characterized by its: (i) computational resource requirement in terms of CPU cycles,    185
denoted by $\gamma_{u_{m,n}}$, (ii) the amount of data uploaded to MEC, denoted by $up_{u_{m,n}}$, (iii) the amount of    186
data that must be downloaded by the user from MEC at the end of execution on the MEC, denoted    187
by $dw_{u_{m,n}}$, and (iv) finally, the maximum tolerated delay according to the Quality of Service (QoS)    188
required by the application, denoted by $t_{u_{m,n}}$.    189

As considered in previous works [11, 19], we distinguish between two categories of computation    190
offloading applications: (1) *static offloading decision task* and (2) *dynamic offloading decision task*. In    191
the first category, the application is partitioned in advance at the design time between: a local part    192
(task) that should always be executed in the mobile and a remote part (task) that should always    193

Table 1. Notation

| Symbol | Description |
|---|---|
| $\gamma_{u_{m,n}}$ | the computational resource required by the task of the user $u_{m,n}$. |
| $\lambda_m$ | the Lagrangian multiplier of the subproblem m. |
| $c_k$ | the computing resource allocation on cloudlet k. |
| $dw_{u_{m,n}}$ | the amount of data to download by the user $u_{m,n}$ from the MEC. |
| $E^l_{u_{m,n}}$ | the local energy consumption for user $u_{m,n}$. |
| $E^e_{u_{m,n},k}$ | the remote energy consumption for user $u_{m,n}$ in cloudlet k. |
| $F_k$ | the computing capacity of the cloudlet k. |
| $f_{u_{m,n}}$ | the local computing capacity of the user $u_{m,n}$. |
| K | the number of cloudlets available in the network. |
| M | the number of APs in network. |
| $N_m$ | the number of users associated to the AP m. |
| $P^{tx/rx}_{u_{m,n}}$ | power consumption when the Wi-Fi interface is transforming or receiving data. |
| $P^{Idle}_{u_{m,n}}$ | power consumption when the Wi-Fi interface is not transforming or receiving data. |
| $r_{u_{m,n}}$ | the task of the user $u_{m,n}$. |
| $t_{u_{m,n}}$ | the maximum tolerated delay according the QoS of the task of the user $u_{m,n}$. |
| $T^t_{u_{m,n},k}$ | the communication time when user $u_{m,n}$ offload to cloudlet k. |
| $T^l_{u_{m,n}}$ | the local processing time for user $u_{m,n}$. |
| $T^e_{u_{m,n},k}$ | the remote processing time for user $u_{m,n}$ in cloudlet k. |
| $u_{m,n}$ | the $n^{th}$ user for the $m^{th}$ AP. |
| $up_{u_{m,n}}$ | the amount of data uploaded to the MEC from the user $u_{m,n}$. |
| $W_m$ | the wireless data rate at the AP m. |
| $w_{u_{m,n}}$ | the allocated data rate for the user $u_{m,n}$. |
| $x_{u_{m,n},k}$ | the offloading decision variable for the task of user $u_{m,n}$ in the cloudlet k. |
| $y_{u_{m,n}}$ | the category belong the tasks (1 static offloading decision task and 0 otherwise). |
| $\mathcal{Z}^l_{u_{m,n}}$ | the local offloading cost for user $u_{m,n}$. |
| $\mathcal{Z}^e_{u_{m,n},k}$ | the remote offloading cost for user $u_{m,n}$ on cloudlet k. |

**Q2**

be executed remotely. As illustrated in Figure 3(a), the task's source code is already in the remote server, so the mobile device needs to transmit only the input data to the remote server. A typical example of static offloading decision task is the FLUID application on Android [16] that is used for particle simulations. The thin client side of FLUID is executed on the mobile, while the server part must be performed remotely in MEC, because it requires high-performance GPU computing processors that are not commonly available in mobile devices.

In the second category, the application needs to be partitioned at runtime accordingly with the network and MEC resource availability. Basically, the mobile terminal needs to decide if it is useful to execute the task on the mobile or to offload all or part of the task. In this case, as illustrated in Figure 3(b), the mobile device must transmit its source code and the input data when the task is offloaded. An example of this kind of application is the Linpack benchmarks [32] for Android, which aims to measure the performances of Android devices. This application can be either totally executed on the mobile terminal or partially offloaded to a cloudlet.

To simplify the analysis, we model in this article both *static offloading decision tasks* and *dynamic offloading decision tasks* as tasks with an offloading computation ratio noted as $a_{u_{m,n}}$. Basically, $a_{u_{m,n}}$ denotes the computation ratio of the application that should be executed remotely. To
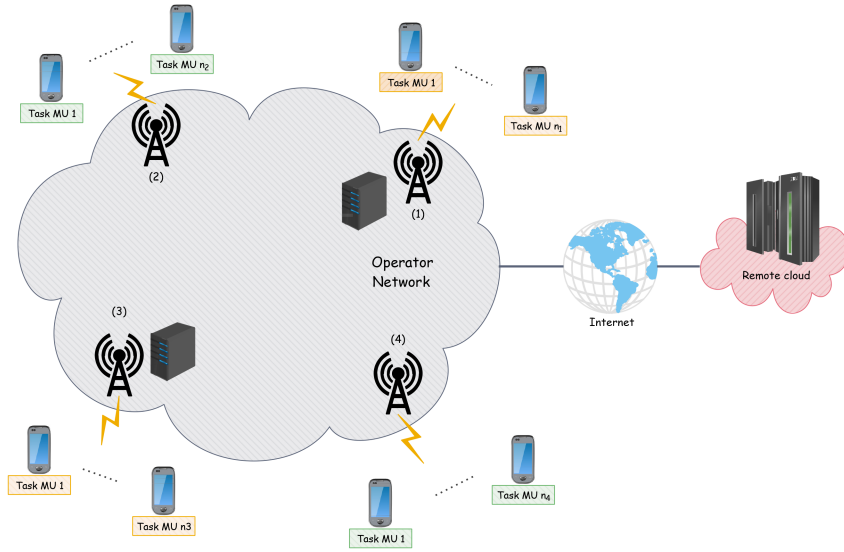
Fig. 2. Example of a multi-user (MU), multi-cloudlet MEC environment with 4 APs and 2 cloudlets ($M = 4$, $K = 2$).



(a) Static offloading decision task           (b) Dynamic offloading decision task
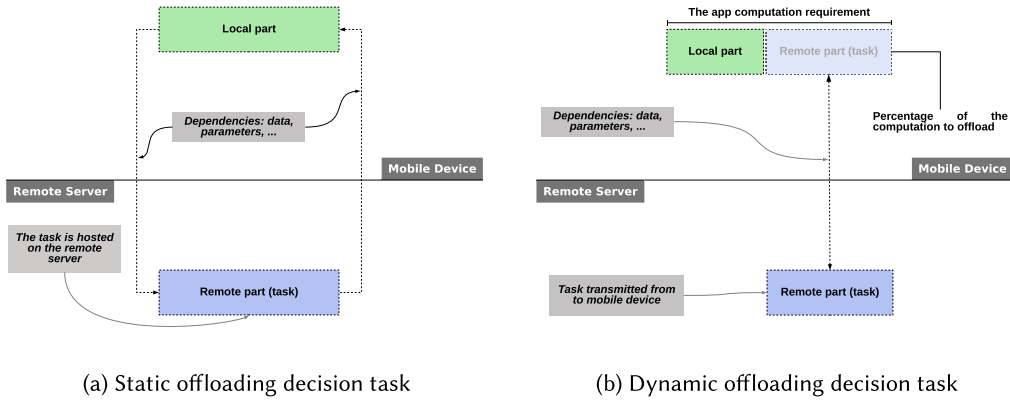
Fig. 3. Illustration of the two categories of tasks.

distinguish between *static offloading decision task*s and *dynamic offloading decision tasks*, we con-  210
sider that for *static offloading decision tasks* the offloading ratio is always equal to 1, which means  211
that the tasks that belong to that category are always offloaded. The local part is considered equal  212
to zero, because it will not affect the performance of the system, since it should always be executed  213
at the terminal. However, for *dynamic offloading decision tasks,* this offloading computation ratio  214
can take any value between 0 and 1 (i.e., $a_{u_{m,n}} \in [0, 1]$). In addition, for simplicity (but without the  215
loss of generality), we also assume that when a task is offloaded with a given offloading computa-  216
tion ratio, then the amount of data that should be transmitted are also proportional to that ratio.  217
However, the output of the task, noted as $dw_{u_{m,n}}$, does not change whatever the value of $a_{u_{m,n}}$.  218

To indicate which category the application of user $u_{m,n}$ belongs, we introduce the binary vari-  219
able $y_{u_{m,n}}$, which is equal to 1 for *static offloading decision task*s and 0 for *dynamic offloading*  220
*decision tasks.*  221

222   Finally, due to hardware and software constraints required by the task, we assume that some
223   cloudlets cannot perform some tasks. In this case, we define a second binary variable, $g_{u_{m,n},k}$, to
224   indicate if the cloudlet $k$ can perform the task. Thus, $g_{u_{m,n},k}$ is equal to 1 if the $k^{th}$ cloudlet can
225   execute the task, 0 otherwise.

### 3.2   Communication Model

227   Let $G = (V, E)$ be a weighted graph, where $V$ ($\mathcal{M} \cup \mathcal{K}$) is a finite set of vertices corresponding to
228   the sets of APs and cloudlets and $E$ is a set of connections (edges) denoting a possible communica-
229   tion between any two vertices. We also consider a weight of the edge, noted as $e_{i,j}$, that represents
230   the network delay between each two vertices $i$ and $j$. Thus, using the Dijkstra algorithm, we can
231   compute a delay matrix, noted as $\mathcal{D}_{m,k}$, that represents the delay between the $m^{th}$ AP and the $k^{th}$
232   cloudlet. According to the last considerations, we can estimate the bandwidth allocated to each
233   user as follows:

$$w_{u_{m,n}} = \frac{W_m(\pi_m)}{\pi_m}, \tag{1}$$

234   where $\pi_m$ is the **offloading capacity** that represents the number of users that offload their ap-
235   plications to MEC, as defined in Definition 5.1. $W_m(\pi_m)$ is the bandwidth shared by the $m^{th}$ AP
236   between the $\pi_m$ users associated with it. To estimate this bandwidth, we consider the Bianchi
237   model [1, 20].
238   According to the last assumptions, we can compute the communication time of any task. This
239   time, noted as $T^t_{u_{m,n},k}$, is composed of time needed to upload the data from the user terminal to the
240   cloudlet, plus the time needed to download the results from the cloudlet once the task is completed.
241   In this case, $T^t_{u_{m,n},k}$ can be written as:

$$T^t_{u_{m,n},k} = a_{u_{m,n}} \cdot \frac{up_{u_{m,n}}}{w_{u_{m,n}}} + \mathcal{D}_{m,k} + \frac{dw_{u_{m,n}}}{w_{u_{m,n}}} + \mathcal{D}_{m,k}$$
$$= \frac{a_{u_{m,n}} \cdot up_{u_{m,n}} + dw_{u_{m,n}}}{w_{u_{m,n}}} + 2\mathcal{D}_{m,k}. \tag{2}$$

### 3.3   Computation Processing Model

243   a) *Local processing:* We assume that the user's device has a local computational capability
244     of $f_{u_{m,n}}$ used for the task computation. When a user offloads a percentage $a_{u_{m,n}}$ of its
245     computation, the remaining part must be performed locally. So, the local processing time
246     for the local part can be estimated by:

$$T^l_{u_{m,n}} = (1 - a_{u_{m,n}}) \cdot \frac{\gamma_{u_{m,n}}}{f_{u_{m,n}}}. \tag{3}$$

247     From the above equation, we can notice that the local processing time of a Static offload de-
248     cision task is equal to zero ($T^l_{u_{m,n}} = 0$), since $a_{u_{m,n}}$ for static offload decision task is always
249     equal to 1.
250   b) *Remote processing:* For remote computation, we consider that each cloudlet has a compu-
251     tational capability of $F_k$. Every cloudlet allocates some of its computational resource to
252     perform the offloaded tasks. So, the remote execution time of the user's $u_{m,n}$ task can be
253     estimated as the ratio of the offloaded computational resources and the allocated compu-
254     tation resource, as illustrated by the following equation:

$$T^e_{u_{m,n},k} = a_{u_{m,n}} \cdot \frac{\gamma_{u_{m,n}}}{c_k}, \tag{4}$$

where $c_k$ is the amount of computational resource allocated to perform the task at the $k^{th}$   255
cloudlet. To simplify our model, we consider that the allocated resource at each cloudlet is   256
fixed and does not change during the computation [4, 5, 12].   257

### 3.4  Offloading Cost Model   258

We define the offloading cost as a combination of the energy consumption and the execution time   259
of the application. In the following, we present both local offloading cost and remote offloading   260
cost.   261

a) *Local offloading cost:* According to the forgoing considerations, the local offloading cost is   262
expressed as a combination of the energy consumption and the execution time of the tasks,   263
as in the following:   264

$$\mathcal{Z}_{u_{m,n}}^{l} = \beta_{u_{m,n}} \cdot E_{u_{m,n}}^{l} + (1 - \beta_{u_{m,n}}) \cdot T_{u_{m,n}}^{l}, \tag{5}$$

where $E_{u_{m,n}}^{l}$ and $T_{u_{m,n}}^{l}$ are, respectively, the total amount of energy and processing time of   265
the local part of the application of the user $u_{m,n}$. $\beta_{u_{m,n}}$ denotes the weighting parameter of   266
execution time and energy consumption of the user's offloading decision. When the battery   267
of the user's device is at a low state, and the user needs to reduce the energy consumption,   268
then it can set $\beta_{u_{m,n}} = 1$. However, when a delay-sensitive task is running, the user can set   269
$\beta_{u_{m,n}} = 0$ to give more priority to the execution time. To get a more flexible cost model,   270
we allow a multi-criteria offloading policy by considering energy consumption, execution   271
time, or a combination of both of them.   272

Considering the energy-consumption model, we use in this article the model proposed   273
in References [3, 23]. Using this model, we can compute $E_{u_{m,n}}^{l}$ as following:   274

$$E_{u_{m,n}}^{l} = \kappa \cdot (f_{u_{m,n}})^3 \cdot T_{u_{m,n}}^{l} \tag{6}$$

where $\kappa$ is the effective switched capacitance, which depends on the chip architecture and   275
is used to adjust the processor frequency. In the following, we set $\kappa = 10^{-9}$ as shown in   276
[3, 23].   277

b) *Remote offloading costs:* The total amount of energy consumed by the user's terminal to   278
perform the task remotely is equal to the energy used when the device turns the radio in   279
the transmission mode to send the data to the remote server, plus the energy used when   280
the device turns the radio in idle mode to wait the task completion plus the energy used   281
by the device when it turn again the radio in the reception mode to receive the result data   282
from the remote server. This consumed energy can be expressed as follows:   283

$$E_{u_{m,n},k}^{e} = P_{u_{m,n}}^{tx/rx} \cdot \left(T_{u_{m,n},k}^{t} - 2\mathcal{D}_{m,k}\right) + P_{u_{m,n}}^{idle} \cdot \left(T_{u_{m,n},k}^{e} + 2\mathcal{D}_{m,k}\right), \tag{7}$$

where $P_{u_{m,n}}^{tx/rx}$ is the power consumption when the radio interface is set to transmission or   284
reception mode, and $P_{u_{m,n}}^{idle}$ is the power consumption in the case when the radio interface   285
is set to idle mode [3, 23, 29].   286

Finally, we can define the remote offloading costs as follows:   287

$$\mathcal{Z}_{u_{m,n},k}^{e} = \beta_{u_{m,n}} \cdot E_{u_{m,n},k}^{e} + (1 - \beta_{u_{m,n}}) \cdot \left(T_{u_{m,n},k}^{t} + T_{u_{m,n},k}^{e}\right) \tag{8}$$

## 4 MULTI-USER, MULTI-CLOUDLET COMPUTATION OFFLOADING PROBLEM FORMULATION AND DECOMPOSITION

To propose an efficient offloading policy, we formulate the problem as an optimization problem. Then, we use Lagrangian relaxation to decompose the problem into subproblems and solve each one separately.

### 4.1 Problem Formulation

As introduced earlier, the objective of this article is to propose an efficient offloading policy. It decides which users should offload their tasks, determines the amount of computation to offload, and selects a cloudlet for each user, while minimizing the total offloading cost. Let us denote $x_{u_{m,n},k}$ to the offloading decision for the task of the user $u_{m,n}$ on the cloudlet $k$, which means that $x_{u_{m,n},k} = 1$ if the user $u_{m,n}$ offloads its task to the cloudlet $k$, 0 otherwise. Given the system description and according to the QoS and cloudlets' resource-capability constraints, the can be formulated as follows:

$$\textbf{Minimize} \sum_{m}^{M} \sum_{n}^{N_m} \mathcal{Z}_{u_{m,n}}$$

**Subject to:**

$$C1 : \sum_{k=1}^{K} x_{u_{m,n},k} \leq 1, \forall m \in \mathcal{M}, u_{m,n} \in \mathcal{N}_m$$

$$C2 : y_{u_{m,n}} - \sum_{k=1}^{K} x_{u_{m,n},k} \leq 0, \forall m \in \mathcal{M}, u_{m,n} \in \mathcal{N}_m$$

$$C3 : T_{u_{m,n}} \leq t_{u_{m,n}}, \forall m \in \mathcal{M}, u_{m,n} \in \mathcal{N}_m$$

$$C4 : x_{u_{m,n},k} \leq g_{u_{m,n},k}, \forall m \in \mathcal{M}, u_{m,n} \in \mathcal{N}_m, k \in \mathcal{K}$$

$$C5 : \sum_{m}^{M} \left( \sum_{n}^{N_m} x_{u_{m,n},k} \cdot c_k \right) \leq F_k, \forall k \in \mathcal{K}$$

$$C6 : x_{u_{m,n},k} \in \{0,1\}, \forall m \in \mathcal{M}, u_{m,n} \in \mathcal{N}_m, k \in \mathcal{K}$$

$$C7 : a_{u_{m,n}} \in [0,1], \forall m \in \mathcal{M}, u_{m,n} \in \mathcal{N}_m$$

$$C8 : a_{u_{m,n}} \geq y_{u_{m,n}}, \forall m \in \mathcal{M}, u_{m,n} \in \mathcal{N}_m, \tag{9}$$

where $\mathcal{Z}_{u_{m,n}}$ is the offloading cost of the user's $u_{m,n}$ task. $\mathcal{Z}_{u_{m,n}}$ can be expressed by the following formula:

$$\mathcal{Z}_{u_{m,n}} = \mathcal{Z}_{u_{m,n}}^{l} + \sum_{k=1}^{K} x_{u_{m,n},k} \cdot \mathcal{Z}_{u_{m,n},k}^{e} \tag{10}$$

As indicated in the problem formulation, our objective is to minimize the total offloading cost of the users of the network. The first constraint ($C1$) ensures that each task is assigned to one cloudlet at most. Constraints ($C2$) guarantee that any static offloading decision task must be assigned to exactly one cloudlet and a dynamic offloading decision task may be assigned to one cloudlet at most. The next constraint ($C3$) shows that the QoS required by the task, in terms of completion time, must be less than a given threshold. The threshold is obtained based on the characteristics of the mobile application [5, 7]. For example, for an interactive application, the user's perception—the duration of the submission of the task until receiving the response—is a well-used technique to

determine the threshold [8]. Completion time can be expressed as follows: 311

$$T_{u_{m,n}} = T^l_{u_{m,n}} + \sum_{k=1}^{K} x_{u_{m,n},k} \cdot \left( T^t_{u_{m,n},k} + T^e_{u_{m,n},k} \right) \tag{11}$$

The next constraint ($C4$) ensures that every offloaded task must be performed by a cloudlet 312
that meets the hardware and software required by the task. The constraint $C5$ shows that it is 313
not possible to exceed the computing capacity of each cloudlet. Constraint $C6$ ensures that any 314
decision variable is a binary variable. 315
Finally, constraints ($C7$) indicate that the ratio of the offloaded computation must be a real value 316
between 0 and 1. And, ($C8$) ensures that each static offloading decision task must always be entirely 317
offloaded. 318

THEOREM 4.1. *Equation (9) is a Non-Linear Mixed Binary Problem (NLMBP) with an exponential* 319
*function and constraints. It is an NP-hard problem.* 320

PROOF. Let us consider a special case where the same number of users are associated to each AP 321
and all tasks are static offloading decisions. So, all the tasks must be offloaded to the cloudlets and 322
the bandwidth allocated to each user is known in advance. Thus, the special case is a Linear Binary 323
Integer Problem (LBIP). In fact, this special case can be easily reduced to the General Assignment 324
Problem (GAP) with assignment restrictions, which is NP-hard, as shown in Reference [17]. Since 325
the special case is NP-hard, Equation (9) is also NP-hard. □ 326

Considering the NP-hardness of the problem, it is difficult to achieve an optimal solution. Next, 327
we propose a simplification version of Equation (9) using Lagrangian relaxation and decomposition 328
approaches. 329

## 4.2 Problem Decomposition 330

To solve the above problem, we need a decomposition approach. Decomposing a complex opti- 331
mization problem consists of breaking it up into smaller ones, called subproblems, and solving 332
each of the smaller ones separately. Unfortunately, the constraint $C5$ is considered a complicat- 333
ing constraint [9, 31], since it involves the local variables of more than one subproblem. Conse- 334
quently, the decomposition of Equation (9) does not work in one step and the subproblem cannot 335
be solved independently. For these kinds of complex problems, there are advanced decomposition 336
techniques that solve the problem by iteratively solving a sequence of subproblems. In this article, 337
we use one of the most popular decomposition techniques, Lagrangian relaxation [9, 31]. The idea 338
of Lagrangian relaxation comes up in the context of using Lagrangian multipliers to decompose 339
the problem; thus, we introduce the Lagrangian multipliers $\lambda = [\lambda_k, k \in \mathcal{K}]^T$ on the constraint $C5$, 340
where $\lambda_k$ denotes the price of all the tasks performed by the $k^{th}$ cloudlet. $X$ and $A$ are the set of 341
the offloading decision variables and the set of the offloading ratio, respectively. The Lagrangian 342
function is given by: 343

$$\begin{aligned} L(X, A, \lambda) &= \sum_{m}^{M} \sum_{n}^{N_m} \mathcal{Z}_{u_{m,n}} + \sum_{k}^{K} \lambda_k \sum_{m}^{M} \sum_{n}^{N_m} (x_{u_{m,n},k} \cdot c_k - F_k), \\ &= \sum_{m}^{M} \sum_{n}^{N_m} \mathcal{Z}_{u_{m,n}} + \sum_{m}^{M} \sum_{n}^{N_m} \sum_{k}^{K} \lambda_k x_{u_{m,n},k} \cdot c_k - \sum_{k}^{K} F_k, \\ &= \sum_{m}^{M} \sum_{n}^{N_m} \left( \mathcal{Z}_{u_{m,n}} + \sum_{k}^{K} \lambda_k x_{u_{m,n},k} \cdot c_k \right) - \sum_{k}^{K} \lambda_k F_k. \end{aligned}$$
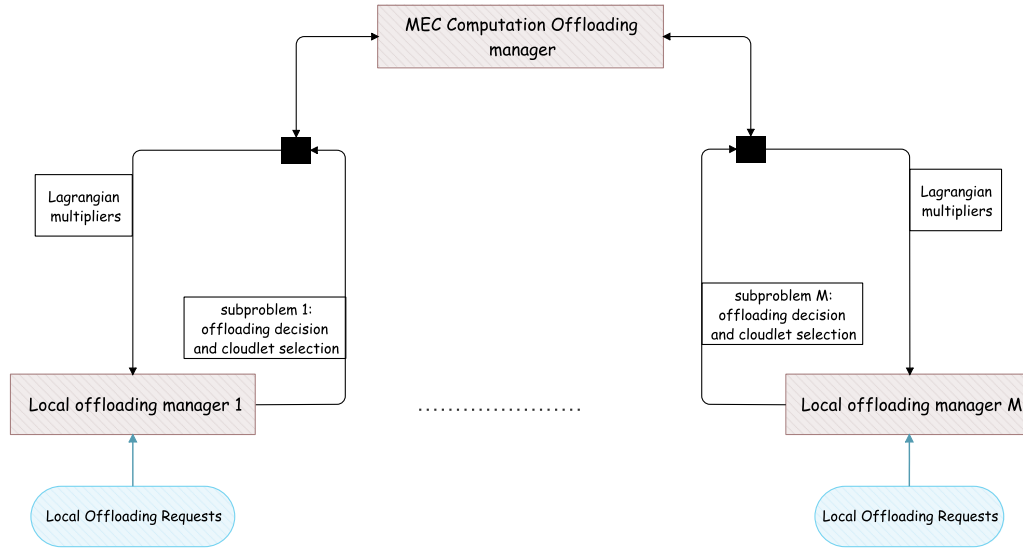
Fig. 4. DM2-ECOP architecture overview.

344  The Lagrangian dual problem for primal Equation (9) is then given by:

$$\max_{\lambda} \min_{X,A} L(X, A, \lambda).$$

345  We can see that the Lagrangian dual problem is separable into two levels: Level 1 is the inner
346  minimizing, which consists of $M$ subproblems, each one concerning only one AP. Level 2 is outer
347  maximization, which is the master problem that considers the global variables and constraint $C5$.
348  Focusing on this observation, we introduce a new offloading policy named Distributed Multi-
349  user Multi-cloudlet Efficient Computation Offloading Policy (DM2-ECOP). In the following, we
350  describe the proposed computation offloading policy.

## 5  DM2-ECOP: DISTRIBUTED MULTI-USER, MULTI-CLOUDLET EFFICIENT COMPUTATION OFFLOADING POLICY

353  As introduced in the last section, the Lagrangian dual problem is decomposable into $M$ subprob-
354  lems. Each subproblem tries to find an optimal offloading decision and cloudlet selection to the
355  users associated with one AP. Considering this characteristic, we design a new offloading policy,
356  DM2-ECOP. As shown in Figure 4, DM2-ECOP has two levels of decision. The local offloading
357  manager is responsible for the offloading decision and cloudlet selection of an AP; it solves the
358  associated subproblem, then sends the offloading decision and the cloudlet selection to the cen-
359  tralized decision level. The MEC computation offloading manager receives the solution of all
360  subproblems, then it ensures that the obtained offloading solution is feasible and respects all the
361  constraints. After, it updates the Lagrangian multipliers and transmits the new values to every
362  local offloading manager to improve the local solutions.

### 5.1  Local Offloading Manager: A Greedy Best Cloudlet Selection First Heuristic(GBC-SFH)

365  The local offloading manager tries to solve the subproblem of one AP to decide which user can
366  offload. Then, it selects the appropriate cloudlet to perform the task of each user. According to
367  the previous considerations, we can formulate the subproblem of a local offloading manager as

follows: 368

$$\textbf{Minimize} \sum_{n}^{N_m} \left( \mathcal{Z}_{u_{m,n}} + \sum_{k}^{K} \lambda_k x_{u_{m,n},k} \cdot c_k \right)$$

**Subject to:**

constraints C1 − C4 and C6−C8. (12)

To solve the subproblem, we need to know the bandwidth allocated to each user. Unfortunately, 369
this bandwidth depends on the number of users that offload their tasks. Therefore, we need to know 370
the bandwidth allocation to decide whether a user should offload its task or not. To overcome this 371
dependency problem, we use a branching heuristic. The key idea is that for any AP $m$, the number 372
of users that can offload their tasks is between an upper and a lower bound. The lower bound 373
corresponds to the number of users with static offloading decision tasks, and the upper bound 374
corresponds to the maximum number of users associated with the AP $m$, $N_m$. 375

*Definition 5.1.* The **offloading capacity** of the AP $m$ is defined as the number of tasks that have 376
been accepted being performed by the MEC environment. We note it by $\pi_m$, and it is given by: 377

$$\pi_m = \sum_{n}^{N_m} \sum_{k}^{K} x_{u_{m,n},k}.$$

The strategy of solving the subproblem is very simple: It consists of finding the optimal $\pi_m$ that 378
gives the minimal offloading cost. We add the constraint $C9$ to the subproblem (12): 379

$$C9 : \sum_{n}^{N_m} \sum_{k}^{K} x_{u_{m,n},k} = \pi_m. \tag{13}$$

To achieve a good and fast offloading decision, the local offloading manager uses greedy heuris- 380
tics to solve the subproblem (12). The Greedy Best Cloudlet Selection First Heuristic (GBC-SFH) 381
offers heuristics to determine which users offload, to determine the computation to offload, and 382
selects the cloudlet to perform each offloaded task. GBC-SFH iterates all possible values of the 383
offloading capacity, $\pi_m$, in an increasing order, as illustrated in Algorithm 1. In brief, the idea is to 384
find the best cloudlet selection for all static offloading tasks in the first step by minimizing the La- 385
grangian cost $\mathcal{Z}_{u_{m,n},k}^e + \lambda_k c_k$ under the constraints $C1 − C4$ and $C6 − C8$. So, each static offloading 386
task is offloaded to the cloudlet that minimizes the Lagrangian cost. 387

For each dynamic offloading decision task, GBC-SFH tries to select the best cloudlet and compute 388
the optimal ratio $a_{u_{m,n}}$ of the computation to offload. According to the resource availability, GBC- 389
SFH can offload the task or perform it locally, $a_{u_{m,n}} = 0$, by the user's device. Since the wireless 390
bandwidth at the AP may not be enough to offload all the dynamic offloading decision tasks, we 391
need to define an order to determine which dynamic offloading decision task is preferred for the 392
offloading. To this purpose, we define an offloading priority for each task according to the following 393
formula: 394

$$\xi_{u_{m,n}} = \mathcal{Z}_{u_{m,n}}^l - \min_{k \in \mathcal{K}}(\mathcal{Z}_{u_{m,n},k}^e); \quad \text{under } a_{u_{m,n}} = 1.$$

Here, the offloading priority is the local cost minus the cost when all the computation is offloaded 395
to the best cloudlet. The idea is that where $\xi_{u_{m,n}}$ is going higher, the user $u_{m,n}$ is more preferred to 396
offload its task. Unlike the static offloading decision tasks, for dynamic offloading decision tasks, 397
we need to compute the computation to offload ($a_{u_{m,n}}$). To this end, GBC-SFH uses a two-step 398
method. In the first step, it selects the best cloudlet to offload the computation of the user $u_{m,n}$. At 399
this step, GBC-SFH chooses the cloudlet that minimizes the Lagrangian cost $\mathcal{Z}_{u_{m,n},k}^e + \lambda_k c_k$ under 400

401  the constraints $C1 - C4$, $C6 - C7$, and $a_{u_{m,n}} = 1$. After the selection of the best cloudlet, GBC-SFH
402  computes the optimal value of $a_{u_{m,n}}$ for the current user. The optimal value of $a_{u_{m,n}}$ is the solution
403  to the following problem:

$$min\left(\mathcal{Z}^e_{u_{m,n},k} + \mathcal{Z}^l_{u_{m,n}}\right)$$

**Subject to:** constraint $C7$. (14)

404  Equation (14) is a simple problem with one variable. Its optimal solution can be achieved by deriv-
405  ative sign rules [28]. Theorem 5.2 shows when the minimum of this problem will be achieved.

406  THEOREM 5.2. *Let us define $\psi_{u_{m,n}}$ and $\mu_{u_{m,n}}$ the upload data-computing ratio of the dynamic of-*
407  *floading decision task, and the local-remote offloading cost ratio of the user $u_{m,n}$, respectively. They*
408  *are given as follows:*

$$\psi_{u_{m,n}} = \frac{up_{u_{m,n}}}{\gamma_{u_{m,n}}},$$

$$\mu_{u_{m,n}} = \frac{w_{u_{m,n}} \cdot [\kappa \cdot f^3_{u_{m,n}} \cdot c_k \cdot \beta_{u_{m,n}} + (1 - \beta_{u_{m,n}}) \cdot (c_k - f_{u_{m,n}}) - \beta_{u_{m,n}} \cdot P^{idle}_{u_{m,n}} \cdot f_{u_{m,n}}]}{c_k \cdot f_{u_{m,n}} \cdot (P^{tx/rx}_{u_{m,n}} \cdot \beta_{u_{m,n}} + 1 - \beta_{u_{m,n}})}$$

409  *The minimum of Equation (14) is achieved when:*

410  • $a_{u_{m,n}} = 1$, *if and only if:*

$$\psi_{u_{m,n}} < \mu_{u_{m,n}}$$

411  • $a_{u_{m,n}} = 0$, *if and only if:*

$$\psi_{u_{m,n}} > \mu_{u_{m,n}}$$

412  • *the problem is constant, if and only if:*

$$\psi_{u_{m,n}} = \mu_{u_{m,n}}.$$

413  PROOF. The proof of Theorem 5.2 is detailed in the appendix. □

414  Using Theorem 5.2, we have three possible scenarios for dynamic offloading decision tasks:
415  (1) when $a_{u_{m,n}} = 1$, the whole computation must be offloaded; (2) when $a_{u_{m,n}} = 0$, there is no
416  offloading; and 3) when the Equation (14) is constant, all possible values of $a_{u_{m,n}}$ give the same
417  performance. As the computation offloading is a solution to improve the performance, we choose
418  to not offload, $a_{u_{m,n}} = 0$ when this case occurs. Once the number of the offloading tasks is equal
419  to the current offloading capacity ($\pi_m$), the remaining tasks are assigned to be performed locally
420  by the user's device.
421  Consequently, in the worst case, the GBC-SFH iterates $N_m$ in the outer loop when there is
422  no static offloading decision task, which means a complexity of $N_m \cdot log(N_m)$ to sort the tasks
423  and $N_m \cdot K$ to assign to task at the inner loop. Thus, the maximum total number of iterations is
424  $N_m^2 \cdot K + N_m^2 \cdot log(N_m)$. Therefore, the complexity of GBC-SFH is $O(N_m^2 \cdot log(N_m))$, which is fast
425  especially when the number of users associated to each AP is small [14, 21, 33] ($\leq$20).

### 5.2 MEC Computation Offloading Manager

427  The outer level of the Lagrangian dual problem is the master problem. It ensures a feasible offload-
428  ing solution of the primal Equation (9). Finding the optimal solution of the Lagrangian dual prob-
429  lem requires an exhaustive search of all solutions' space and Lagrangian multiplier values, which
430  is a difficult task in general [9]. Consequently, we need to adopt a faster approach. In this work, we
431  use the Subgradient-based heuristic [31]. The proposed heuristic used in the MEC computation
432  offloading manager has three steps, as illustrated in Algorithm 2. First, it solves the subproblems

---

**ALGORITHM 1:** *The local offloading manager: GBC-SFH*

---

**Input:**
1: $\Pi_m$: Set of offloading capacity;
**Output:** Output the offloading decision $X$, ratio $\mathcal{A}$, and cost $\mathcal{Z}$;
2: Sort $\Pi_m$ in increasing order;
3: **for** $\pi_m \in \Pi_m$, **do**
4:   allocate bandwidth using Equation (1);
5:   offload each static offloading decision task to the cloudlet $k$ that minimizes $\mathcal{Z}^e_{u_{m,n},k} + \lambda_k c_k$
     under constraints $C1 - C4$ and $C6 - C8$;
6:   $nb_{offloaded\_task}$ = number of static offloading decision tasks;
7:   compute $\xi_{u_{m,n}}$ for every dynamic offloading decision task;
8:   Sort dynamic offloading decision tasks in decreasing order of $\xi_{n_m}$;
9:   **while** $nb_{offloaded\_task} \leq \pi_m$, **do**
10:     Select the cloudlet $k$ that minimizes $\mathcal{Z}^e_{u_{m,n},k} + \lambda_k c_k$ under constraints $C1 - C4$, $C6 - C7$,
        and $a_{u_{m,n}} = 1$;
11:     Compute the optimal value of $a_{u_{m,n}}$ using Theorem 5.2;
12:     **if** $a_{u_{m,n}} == 0$, **then**
13:       there is no offloading. This dynamic Offloading task must be performed locally;
14:     **else**
15:       Offload this dynamic Offloading task to the cloudlet $k$;
16:       $nb_{offloaded\_task} + +$;
17:     **end if**
18:     **if** (there is no more task) and ($nb_{offloaded\_task} < \pi_m$), **then**
19:       break the while-loop. There is no feasible solution for this value of $\pi_m$;
20:     **end if**
21:   **end while**
22:   all the remaining tasks must be performed locally;
23:   update the best offloading cost $\mathcal{Z}$, ratio $\mathcal{A}$ and decision $X$;
24: **end for**

---

in the local offloading manager by the GBC-SFH for the current Lagrangian multipliers $\lambda$. Next, 433
the MEC computation offloading manager checks if they obtained an offloading solution that 434
is not feasible. If so, the Lagrangian Adjustment Heuristic (LAH) will be used the get a feasible 435
solution using a local search. The idea of LAH heuristics is to check if every cloudlet respects 436
the constraint $C5$. When a cloudlet does not respect this constraint, LAH heuristic reassigns some 437
tasks offloaded from this cloudlet to another cloudlet that respects all constraints. 438

At the end, the MEC computation offloading manager updates the Lagrangian multipliers 439
by the following formula: 440

$$\lambda_k(t+1) = \lambda_k(t) + \theta(t) \cdot \left( \sum_m^M \left( \sum_n^{N_m} x_{u_{m,n},k} \cdot c_k \right) - F_k \right), \tag{15}$$

where $\theta(t)$ is the update step. In this work, we use the Held and Karp formula [9, 31] to update 441
this step as follows: 442

$$\theta(t) = \eta(t) \cdot \frac{\mathcal{Z}^* - \mathcal{Z}(t)}{\sum_{k=1}^K (\sum_m^M \sum_n^{N_m} x_{u_{m,n},k} \cdot c_k - F_k)^2}, \tag{16}$$

---

**ALGORITHM 2:** *MEC computation offloading manager*

---

**Input:**
  1: $It_{max}$: maximum number of iterations;
  2: $\varepsilon$: an infinitesimal number;
**Output:** offloading decision and cloudlet selection for all users;
  3: Initialize $\lambda_k$ randomly;
  4: $\mathcal{Z}_{max} = -\infty$;
  5: **while** (t $< It_{max}$ and $\theta(t) > \varepsilon$), **do**
  6:    **for** ($m \in \mathcal{M}$), **do**
  7:      $\mathcal{Z}_m(t)$= get the solution of subproblem $m$ from the local offloading manager $m$;
  8:    **end for**
  9:    $\mathcal{Z}(t) = \sum_{m \in \mathcal{M}} \mathcal{Z}(t)_m - \sum_{k \in \mathcal{K}} \lambda_k \cdot F_k$;
 10:    **if** ($\mathcal{Z}(t) > \mathcal{Z}_{max}$), **then**
 11:      $\mathcal{Z}_{feasible}$ = use Heuristic *LAH* to find a feasible solution;
 12:      **if** ($\mathcal{Z}_{feasible} < \mathcal{Z}^*$), **then**
 13:        $\mathcal{Z}^* = \mathcal{Z}_{feasible}$;
 14:        update the best solution of the primal problem;
 15:      **end if**
 16:      $\mathcal{Z}_{max} = \mathcal{Z}(t)$;
 17:    **end if**
 18:    update the Lagrangian multipliers and $\eta$ using Equations (15), (16), and (17);
 19: **end while**

---

where $\eta(t)$ is a decreasing adaptation parameter $0 < \eta(0) \le 2$, $\mathcal{Z}^*$ is the best obtained solution of Equation (9), and $\mathcal{Z}(t)$ refers to the current solution of the Lagrangian dual problem. $\eta(t)$ can be expressed by the following formula:

$$\eta(t+1) = \begin{cases} \vartheta \cdot \eta(t) & \text{if } \mathcal{Z}(t) \text{ did not increase} \\ \eta(t) & \text{otherwise} \end{cases}. \tag{17}$$

As suggested in References [9, 31], we set the values of $\vartheta = 0.9$ and $\eta(0) = 2$. The master problem repeats these steps until the stop conditions, which are the maximum number of iterations $It_{max}$ and the maximum tolerated error of the update step $\varepsilon$ ($\theta \le \varepsilon$).

## 6 NUMERICAL RESULTS

In this section, we evaluate the performance of DM2-ECOP using the characteristics of realistic system configuration. We use an MEC environment consisting of a metropolitan area, which is composed of 20 APs forming a ring topology. The delay between any two APs is 3*ms* and the delay between every AP and the remote cloud is 100*ms* [14, 33]. We suppose that four cloudlets are equidistantly deployed among the network, i.e., cloudlet 1 is collocated with the AP 1, cloudlets 2 with the AP 6, cloudlet 3 with the AP 11, and cloudlet 4 with the AP 16. To study the performance of our offloading policy, we consider four cloudlet configurations. Table 2 illustrates the list of the cloudlets' configurations considered for our tests. We consider real configurations used by public cloud providers, such as, Amazon Web Services (AWS) and Microsoft Azure [5, 12, 21], to simulate the behavior of DM2-ECOP policy for real-world scenarios.

The wireless bandwidth of each AP is 150Mbps. The bandwidth allocated to each user is estimated using the parameter settings used in the Bianchi model [1]. Similar to Reference [33], we

Table 2.  List of the Cloudlets' Configurations Used for the Tests

| Configuration | Computing capacity $F_k$ and allocation $c_k$ in Giga CPU cycles/s | | | | | | | |
| | cloudlet 1 | | cloudlet 2 | | cloudlet 3 | | cloudlet 4 | |
| | $c_1$ | $F_1$ | $c_2$ | $F_2$ | $c_3$ | $F_3$ | $c_4$ | $F_4$ |
|---|---|---|---|---|---|---|---|---|
| configuration 1 | 10 | 1,000 | 10 | 1,000 | 10 | 1,000 | 10 | 1,000 |
| configuration 2 | 15 | 1,000 | 10 | 1,000 | 15 | 1,000 | 10 | 1,000 |
| configuration 3 | 10 | 500 | 10 | 500 | 10 | 1,500 | 10 | 1,500 |
| configuration 4 | 15 | 500 | 10 | 500 | 15 | 1,500 | 10 | 1,500 |

Table 3.  The Characteristic of the Real-world Applications
Used for Our Tests

| Application | $\gamma u_{m,n}$ (Giga CPU cycles) | $up_{u_{m,n}}$ (Kilobyte) | $dw_{u_{m,n}}$ (Byte) | $t_{u_{m,n}}$ (Second) |
|---|---|---|---|---|
| static offloading decision tasks | | | | |
| FACE | 12.3 | 62 | 60 | 5 |
| SPEECH | 15 | 243 | 50 | 5.1 |
| OBJECT | 44.6 | 73 | 50 | 13 |
| dynamic offloading decision tasks | | | | |
| Linpack | 50 | 10,240 | 120 | 62.5 |
| CPUBENCH | 3.36 | 80 | 80 | 4.21 |
| PI BENCH | 130 | 10,240 | 200 | 163 |

assume that the number of users connected to every AP, $N_m$ is not greater than 20. Precisely, $N_m$ takes values from $\{5, 10, 15, 20\}$. Each user runs one application from Table 3. The first three applications are static offloading decision tasks; the others are dynamic offloading decision tasks [11]. We assume that $P_{u_{m,n}}^{tx/rx} = 10 * P_{u_{m,n}}^{idle}$ and $P_{u_{m,n}}^{idle} = 100mW$, as shown in Reference [3]. The local computing capability of each user was randomly chosen from $f_{u_{m,n}} \in [0.8, 1, 1.2]$ gigacycles.

The performances of DM2-ECOP are compared to two offloading policies from the literature:

- Nearest Cloudlet Offloading (NCO) [14, 33]: in which each AP is associated with the nearest cloudlet. So, all the users connected to this AP offload their tasks to the same cloudlet. When a cloudlet is overloaded, the tasks are migrated to another cloudlet.
- Full Cloud Offloading (FCO) [4, 5]: In this case, the users offload their tasks to the remote cloud. To make sense of the performances comparison of the offloading policies DM2-ECOP, NCO, and FCO, we assume that the computing capacity allocated to perform each offloading task in the remote cloud is 10 gigacycles.
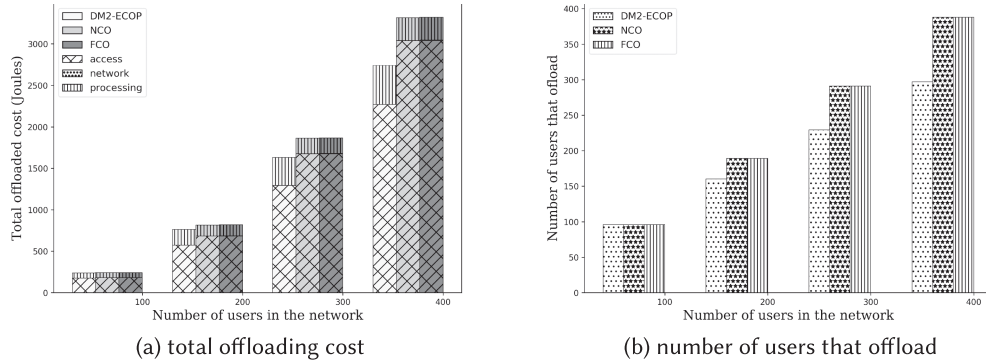
In the following, the default cloudlet configuration is the first configuration (configuration 1), and the density of users at each AP is considered as the same, i.e., the same number of users at each AP. Furthermore, the stop criteria of the MEC computation offloading manager for DM2-ECOP are $It_{max} = 100$ for the maximal number of iterations, and $\varepsilon = 10^{-20}$ for the maximum tolerated error of update steps.

## 6.1 Convergence of DM2-ECOP

To evaluate the performance of DM2-ECOP and its convergence to a feasible solution, Table 4 depicts the required number of iterations to get a feasible solution and the last value of update steps ($\theta$). As expected, the required number of iterations increases as the number of users in the

Table 4. Number of Iterations and Update Step Taken by
DM2-ECOP to Converge to a Feasible Solution

| Number of users | Number of iterations | Update step $\theta$ |
|---|---|---|
| 100 | 15 | 0.0 |
| 200 | 20 | $9.12 \times 10^{-22}$ |
| 300 | 29 | $5.09 \times 10^{-21}$ |
| 400 | 43 | $4.46 \times 10^{-21}$ |



(a) total offloading cost

(b) number of users that offload

Fig. 5. Comparison of offloading policies DM2-ECOP, NCO, and FCO where the cost parameter $\beta = 1$.

network increases. Moreover, the update step converges to the maximum tolerated error $\varepsilon$ with a few number of iterations; this convergence changes while the number of users increases. Thanks to the Held and Karp formula used in our work to update the Lagrangian update step, for the rapid convergence of DM2-ECOP offloading policy.

## 6.2 Offloading Performance Comparison

Figure 5 plots the offloading performances of DM2-ECOP, *NCO*, and *FCO* if we set the cost parameter $\beta = 1$. We also distinguish between the costs related to the network access, network backhaul, and processing. We note that DM2-ECOP reduces the total offloading cost compared to NCO and FCO. More precisely, we can see that the access cost of DM2-ECOP is the lowest, but the processing cost is the highest. This is due to the bandwidth allocation heuristic used by DM2-ECOP, which tries to maximize the bandwidth allocated to each user by minimizing the offloading capacity ($\pi$) of each AP. So, less users can offload their tasks to the MEC with DM2-ECOP compared to *NCO* and *FCO*. However, where the wireless bandwidth is enough to offload all tasks, DM2-ECOP and *NCO* are equivalent, as shown in Table 4 where 100 users are in the network.

To understand the effect of user density on offloading performances, we investigate in Figure 6 the offloading gain of DM2-ECOP compared to NCO under different user density. We consider four scenarios where the topology is divided into two regions, each one containing 10 APs. In Scenario 1, the regions have the same user density. In Scenario 2, user density in Region 1 is twice the user density in Region 2. In Scenario 3, user density in Region 1 is three times the user density in Region 2. Finally, in Scenario 4, user density in Region 1 is four times the user density in Region 2. In Figure 6, we note that the offloading gain of DM2-ECOP compared to NCO goes up when user density goes high. For example, where 200 users are in the network, the gain is 6.1% for Scenario 1 and 10.5% for Scenario 4. This is because the cloudlet selection in *NCO* is static, so this policy needs to migrate some tasks where the cloudlet is overloaded. Consequently, it adds an extra offloading
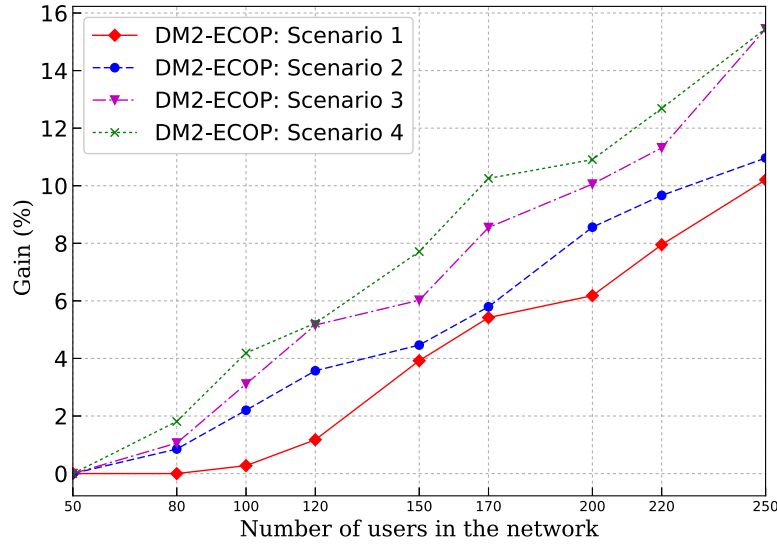
Fig. 6.  Comparison of DM2-ECOP and NCO for different user density in the network.
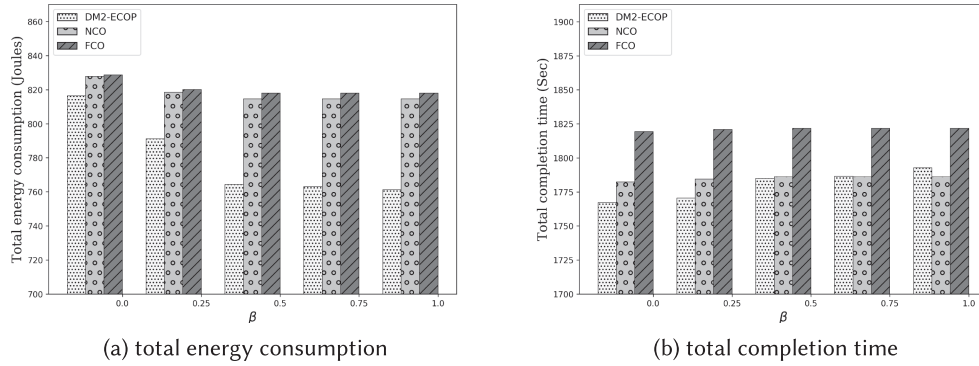


(a) total energy consumption

(b) total completion time

Fig. 7.  Comparison of offloading policies DM2-ECOP, NCO, and FCO over the parameter $\beta$, where 200 users are in the mobile edge computing environment.

cost. However, DM2-ECOP tries to find the best cloudlet selections dynamically at the offloading    508
decision, according to the system and network resource availability.    509

### 6.3    Impact of the Cost Parameter $\beta$ on the Offloading Performance    510

Figure 7 studies the effect of the offloading cost parameter $\beta$ on the performance of the policies    511
DM2-ECOP, *NCO*, and *FCO*. As we can see in Figure 7(a), the energy consumption of DM2-ECOP    512
is better than *NCO* and *FCO* for all possible values of $\beta$. Indeed, even if we set $\beta$ to 0, which    513
means that we give a complete priority to the tasks' completion time, DM2-ECOP obtains better    514
performances. Moreover, when we increase the value of $\beta$, the obtained performances are even    515
better. Consequently, DM2-ECOP achieves better performance, in terms of energy consumption,    516
whatever the offloading cost: energy consumption, completion time, or a combination of energy    517
and time. This is because of the dynamic cloudlet selection adopted in DM2-ECOP.    518

In Figure 7(b), we investigate the effect of $\beta$ on the performance in terms of completion time. We    519
note that the completion time of DM2-ECOP and NCO are better than FCO, because the cloudlets    520

(a) total energy consumption                    (b) total completion time
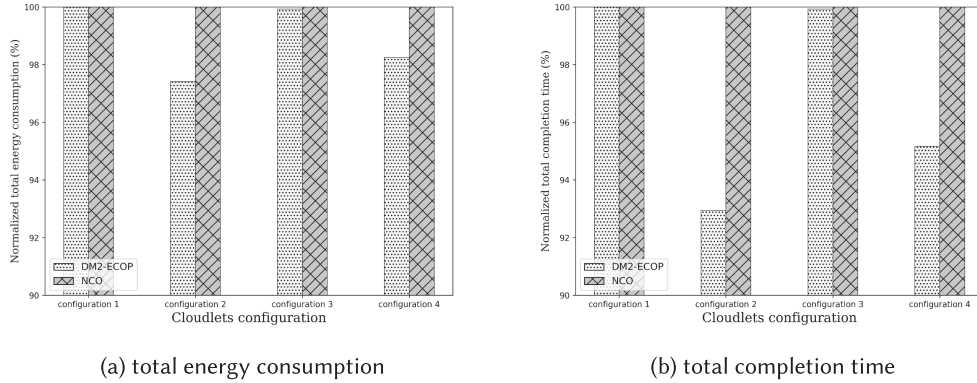
Fig. 8. Offloading policies' performances over different cloudlet configurations, where $\beta = 1$ (the offloading cost is energy consumption) and 200 users are in the MEC environment.

are close to users. Moreover, the completion time of DM2-ECOP is the lowest where $\beta = 0$. However, NCO achieves better completion time than DM2-ECOP where $\beta$ closes to 1. Consequently, NCO has the best performance in terms of completion time. In fact, when $\beta$ is close to 1, the energy consumption becomes more important than the completion time in the expression of the offloading cost. DM2-ECOP reduces the offloading cost by offloading less tasks to the MEC, as shown in Figure 5, to minimize the energy consumed by the wireless access level. As a result, more tasks are executed locally, which increases the completion time.

### 6.4   Impact of the Cloudlets' Configurations

In the following, we study the performance of the offloading policies over heterogeneous cloudlet configurations. In Figure 8, we investigate the performance of the offloading policies DM2-ECOP and *NCO* over four cloudlet configurations presented at the beginning of this section. We observe that the DM2-ECOP and *NCO* are equivalent when the cloudlets have exactly the same configurations, such as Configuration 1. However, where the computing resources allocated to each task are heterogeneous, which corresponds to a more realistic scenario, DM2-ECOP achieves better performance in terms of energy consumption and completion time. This can be explained by the fact that DM2-ECOP offloads to the best cloudlet, but *NCO* offloads to the nearest cloudlet. We also can notice that DM2-ECOP gets better performance for Configuration 2 than Configuration 4, even if the resource allocated for each task are the same in the two configurations. This is because the total computing capacity in Configuration 4 is not homogeneous. Thus, NCO needs to migrate some tasks from Cloudlets 1 and 2 to Cloudlets 3 and 4. To summarize, these results show that DM2-ECOP can achieve a good offloading performance under different cloudlet configurations.

### 6.5   Impact of the Applications Characteristics

As shown previously in our analytical model, the characteristics of the application have a crucial role in the efficiency of the offloading performance. To understand this role, we investigate the impact of the application on the offloading cost and on the amount of the offloaded computation, $a_{u_{m,n}}$.

Figure 9 depicts the performances of the offloading policies for each application under the cloudlets' Configuration 3. We note that for static offloading decision tasks, DM2-ECOP has better energy consumption and completion time followed by *NCO*. Indeed, the static offloading decision tasks must be offloaded, and DM2-ECOP tries to find the best cloudlet selection at the decision

(a) total energy consumption                    (b) total completion time
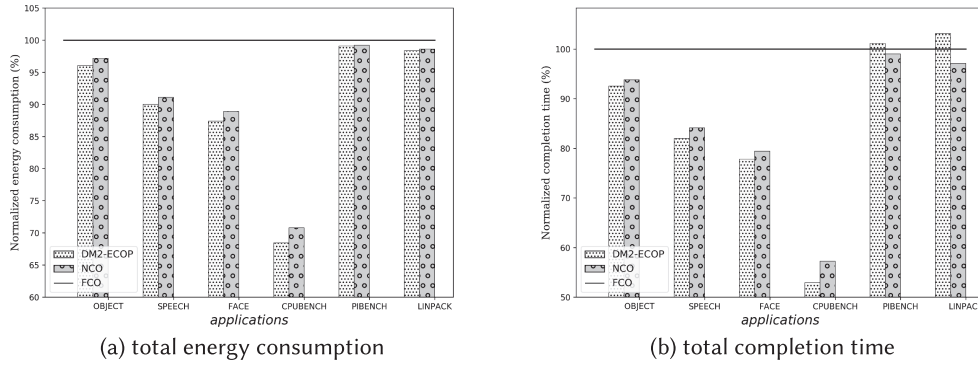
Fig. 9.  Offloading policies' performances for each application for Configuration 3, where $\beta = 1$ (the offloading cost is energy consumption) and 200 users are in the MEC environment.
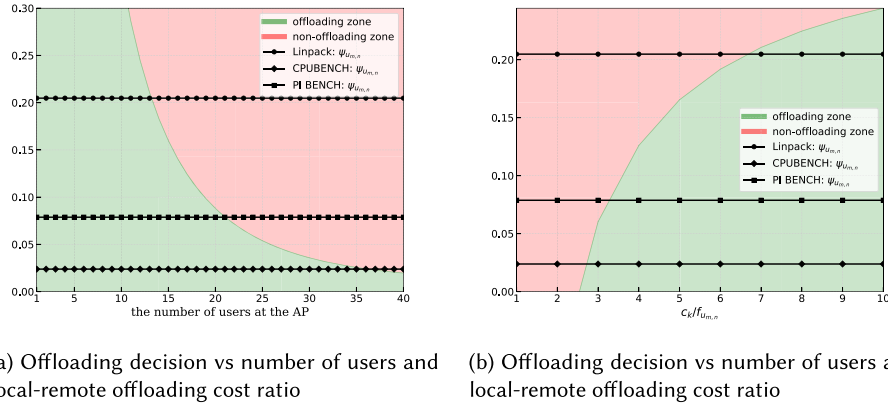


(a) Offloading decision vs number of users and local-remote offloading cost ratio

(b) Offloading decision vs number of users and local-remote offloading cost ratio

Fig. 10.  The effect of the number of users, $w_{u_{m,n}}$, the considered offloading cost $\beta_{u_{m,n}}$, and the allocation computing resource $c_k$ on the application partition decision, where the wireless bandwidth at AP is 150mbps, local cpu power is $f_{u_{m,n}} = 1$ gigacycle.

time. However, FCO offloads to the nearest cloudlet, which induces an additional offloading cost, since some tasks must be migrated to other cloudlets. However, for the dynamic offloading decision tasks, we note that DM2-ECOP has the lowest energy consumption and completion time where the application does not need to upload lots of data to the cloudlet, such as CPUBENCH. However, where the application requires lots of data, e.g., Linpack and PIBENCH, the completion time of *NCO* and *FCO* are better than that of DM2-ECOP. This is because DM2-ECOP tries to perform the application locally when it uploads a lot of data to minimize the access cost of the users.

Finally, to deeply analyze the performance of DM2-ECOP with the dynamic offloading decision tasks, we study the effect of the number of users and the ratio between the remote and local processing capacity on the offloading decision. In Figure 10(a), we plot the effect of the number of users per access point and, thus, the amount of bandwidth allocated to each user on the offloading decision. The red area corresponds to the case where we execute a task only locally (i.e., $a_{u_{m,n}} = 0$), and the green area corresponds to the case where a task is totally offloaded (i.e., $a_{u_{m,n}} = 1$). In addition to the offloading decision, we also plot the upload data-computing ratio $\psi_{u_{m,n}}$ for three dynamic offloading decision tasks, namely Linpack, CPUBENCH, and PI BENCH. As we can see,

the three applications do not behave the same way when we increase the number of users per AP. Indeed, Linpack is the most sensitive application and stops offloading when the number of users is greater than 14. However, CPUBENCH is the less-sensitive application, since it stops offloading when the number of users reaches 37. This is due to the fact that Linpack sends much more data when it offloads compared to CPUBENCH. In Figure 10(b), we investigate the impact of the ratio between the local and remote processing capacity on the offloading decision. As in the last figure, we also plot the upload data-computing ratio $\psi_{u_{m,n}}$ for Linpack, CPUBENCH, and PI BENCH. As we can see, more the remote processing capacity is important compared to the local processing capacity more the decision is to offload task. However, as in the last figure, Linpack is less sensitive to that increase compared to CPUBENCH and PI BENCH. Indeed, since the total amount of data that should be offloaded for Linpack is important, the offloading becomes beneficial only if the remote processing capability is very important in comparison to the local processing capability.

## 7   CONCLUSION

Computation offloading in a multi-user, multi-cloudlet mobile edge computing environment is a challenging issue. In this article, we propose a new computation offloading policy to decide which users should offload and to which cloudlet. First, we formulate the problem as a Non-Linear Mixed Binary Integer Program. Then, we propose an efficient distributed heuristic to solve the problem using the Lagrangian decomposition approach. The proposed heuristic uses a branching algorithm to maximize the bandwidth allocation and minimize the offloading cost.

In addition, compared to other works, our proposal (DM2-ECOP) considers two categories of offloadable tasks: the static offloading decision tasks that must be performed remotely and the dynamic offloading decision tasks that can be performed both locally and remotely. We also add an offloading computation ratio associated with both static and dynamic decision tasks. This ratio denotes the portion of the application that is executed locally in the terminal and the portion of the application that should be offloaded to the cloudlet.

The obtained numerical results show performance improvements in terms of the offloading cost compared to existing offloading policies under different scenarios and cloudlet configurations. Moreover, because we consider that all the tasks have the same priority and they are not sharing resources (same CPU) at the cloudlet, we demonstrate that the best possible value of the offloading computation ratio is either 0 or 1.

For future work, we will consider an adaptive offloading policy, where the offloaded tasks must be determined at runtime. Moreover, in this article, we assume that each mobile is executing only one task at a time. In future work, we propose to explore the case where an application is characterized by a tasks dependency graph. In this case, more than one task can be offloaded at the same time to the remote cloudlet or cloud.

## APPENDIX: PROOF OF THEOREM 5.2

Given the objective function of Equation (14), which is a function with the variable $a_{u_{m,n}}$, we can find the minimum following the derivative sign rules [28]. Let $\mathcal{F}_{u_{m,n}}$ be the derivative of the objective function, $\mathcal{Z}^e_{u_{m,n},k} + \mathcal{Z}^l_{u_{m,n}}$. $\mathcal{F}_{u_{m,n}}$ is given as follows:

$$
\mathcal{F}_{u_{m,n}} = up_{u_{m,n}} \cdot \frac{P^{tx/rx}_{u_{m,n}} \cdot \beta_{u_{m,n}} + 1 - \beta_{u_{m,n}}}{w_{u_{m,n}}}
$$
$$
+ \gamma_{u_{m,n}} \cdot \left( \frac{\beta_{u_{m,n}} \cdot P^{idle}_{u_{m,n}}}{c_k} + \frac{1 - \beta_{u_{m,n}}}{c_k} - \frac{1 - \beta_{u_{m,n}}}{f_{u_{m,n}}} - \kappa \cdot f^2_{u_{m,n}} \cdot \beta_{u_{m,n}} \right).
$$

The derivative function $\mathcal{F}_{u_{m,n}}$ is a constant and does not change when the variable $a_{u_{m,n}}$ does. 606
According to the derivative sign rules [28], we have three cases: 607

**Case 1:** $\mathcal{F}_{u_{m,n}} < 0$, there the objective function of Equation (14) is monotonically decreasing. 608
Consequently, its minimum is achieved at $a_{u_{m,n}} = 1$. This case occurs when: 609

$$\frac{up_{u_{m,n}}}{\gamma_{u_{m,n}}} < \frac{w_{u_{m,n}} \cdot \left[\kappa \cdot f^3_{u_{m,n}} \cdot c_k \cdot \beta_{u_{m,n}} + (1 - \beta_{u_{m,n}}) \cdot (c_k - f_{u_{m,n}}) - \beta_{u_{m,n}} \cdot P^{idle}_{u_{m,n}} \cdot f_{u_{m,n}}\right]}{c_k \cdot f_{u_{m,n}} \cdot \left(P^{tx/rx}_{u_{m,n}} \cdot \beta_{u_{m,n}} + 1 - \beta_{u_{m,n}}\right)}.$$

**Case 2:** $\mathcal{F}_{u_{m,n}} > 0$, there the objective function of Equation (14) is monotonically increasing. 610
Consequently, its minimum is achieved at $a_{u_{m,n}} = 0$. This case occurs when: 611

$$\frac{up_{u_{m,n}}}{\gamma_{u_{m,n}}} > \frac{w_{u_{m,n}} \cdot \left[\kappa \cdot f^3_{u_{m,n}} \cdot c_k \cdot \beta_{u_{m,n}} + (1 - \beta_{u_{m,n}}) \cdot (c_k - f_{u_{m,n}}) - \beta_{u_{m,n}} \cdot P^{idle}_{u_{m,n}} \cdot f_{u_{m,n}}\right]}{c_k \cdot f_{u_{m,n}} \cdot \left(P^{tx/rx}_{u_{m,n}} \cdot \beta_{u_{m,n}} + 1 - \beta_{u_{m,n}}\right)}.$$

**Case 3:** $\mathcal{F}_{u_{m,n}} = 0$, there the objective function of Equation (14) is constant and does not change. 612
So, at the values of $a_{u_{m,n}}$ give the same cost. This case occurs when: 613

$$\frac{up_{u_{m,n}}}{\gamma_{u_{m,n}}} = \frac{w_{u_{m,n}} \cdot \left[\kappa \cdot f^3_{u_{m,n}} \cdot c_k \cdot \beta_{u_{m,n}} + (1 - \beta_{u_{m,n}}) \cdot (c_k - f_{u_{m,n}}) - \beta_{u_{m,n}} \cdot P^{idle}_{u_{m,n}} \cdot f_{u_{m,n}}\right]}{c_k \cdot f_{u_{m,n}} \cdot \left(P^{tx/rx}_{u_{m,n}} \cdot \beta_{u_{m,n}} + 1 - \beta_{u_{m,n}}\right)}.$$

Ending of the proof. □ 614

## REFERENCES

[1] Giuseppe Bianchi. 2000. Performance analysis of the IEEE 802.11 distributed coordination function. *IEEE J. Select. Areas Comm.* 18, 3 (2000), 535–547. 615 616

[2] Arash Bozorgchenani, Daniele Tarchi, and Giovanni Emanuele Corazza. 2017. An energy and delay-efficient partial offloading technique for fog computing architectures. In *Proceedings of the IEEE Global Communications Conference (GLOBECOM'17)*. IEEE, 1–6. 617 618 619

[3] Aaron Carroll, Gernot Heiser. 2010. An analysis of power consumption in a smartphone. In *Proceedings of the USENIX Annual Technical Conference (USENIXATC'10)*, Vol. 14. Boston, MA, 21–21. 620 621

[4] Meng-Hsi Chen, Ben Liang, and Min Dong. 2016. Joint offloading decision and resource allocation for multi-user multi-task mobile cloud. In *Proceedings of the IEEE International Conference on Communications (ICC'16)*. IEEE, 1–6. 622 623

[5] Xu Chen, Lei Jiao, Wenzhong Li, and Xiaoming Fu. 2016. Efficient multi-user computation offloading for mobile-edge cloud computing. *IEEE/ACM Trans. Netw.* 24, 5 (2016), 2795–2808. 624 625

[6] Byung-Gon Chun, Sunghwan Ihm, Petros Maniatis, Mayur Naik, and Ashwin Patti. 2011. Clonecloud: Elastic execution between mobile device and cloud. In *Proceedings of the 6th European Conference on Computer Systems (EuroSys'11)*. ACM, 301–314. 626 627 628

[7] Eduardo Cuervo, Aruna Balasubramanian, Dae-ki Cho, Alec Wolman, Stefan Saroiu, Ranveer Chandra, and Paramvir Bahl. 2010. MAUI: Making smartphones last longer with code offload. In *Proceedings of the 8th International Conference on Mobile Systems, Applications, and Services (MobiSys'10)*. ACM, 49–62. 629 630 631

[8] Debessay Fesehaye, Yunlong Gao, Klara Nahrstedt, and Guijun Wang. 2012. Impact of cloudlets on interactive mobile cloud applications. In *Proceedings of the 16th IEEE International Enterprise Distributed Object Computing Conference (EDOC'12)*. IEEE, 123–132. 632 633 634

[9] Marshall L. Fisher. 2004. The Lagrangian relaxation method for solving integer programming problems. *Manag. Sci.* 50, 12-supplement (2004), 1861–1871. 635 636

[10] Keke Gai, Meikang Qiu, and Hui Zhao. 2018. Energy-aware task assignment for mobile cyber-enabled applications in heterogeneous cloud computing. *J. Parallel and Distrib. Comput.* 111 (2018), 126–135. 637 638

[11] Ying Gao, Wenlu Hu, Kiryong Ha, Brandon Amos, Padmanabhan Pillai, and Mahadev Satyanarayanan. 2015. Are cloudlets necessary? October (2015). 639 640

[12] Songtao Guo, Bin Xiao, Yuanyuan Yang, and Yang Yang. 2016. Energy-efficient dynamic offloading and resource scheduling in mobile cloud computing. In *Proceedings of the 35th IEEE International Conference on Computer Communications (INFOCOM'16)*, Vol. 2016-July. IEEE, 1–9. 641 642 643

[13] Dong Huang, Ping Wang, and Dusit Niyato. 2012. A dynamic offloading algorithm for mobile computing. *IEEE Trans. Wireless Comm.* 11, 6 (2012). IEEE, 1991–1995. 644 645

**Q6**

[14] Mike Jia, Jiannong Cao, and Weifa Liang. 2015. Optimal cloudlet placement and user to cloudlet allocation in wireless metropolitan area networks. *IEEE Trans. Cloud Comput.* 99 (2015).

[15] Mike Jia, Weifa Liang, Zichuan Xu, and Meitian Huang. 2016. Cloudlet load balancing in wireless metropolitan area networks. In *Proceedings of the 35th IEEE International Conference on Computer Communications (INFOCOM'16)*, Vol. 2016-July. IEEE, 1–9.

[16] Doyub Kim, Woojong Koh, Rahul Narain, Kayvon Fatahalian, Adrien Treuille, and James F. O'Brien. 2013. Near-exhaustive precomputation of secondary cloth effects. *ACM Trans. Graph.* 32, 4 (July 2013), 87:1–87:8.

[17] Sven O. Krumke and Clemens Thielen. 2013. The generalized assignment problem with minimum quantities. *Euro. J. Op. Res.* 228, 1 (2013), 46–55.

[18] Grace Lewis, Sebastián Echeverría, Soumya Simanta, Ben Bradshaw, and James Root. 2014. Tactical cloudlets: Moving cloud computing to the edge. In *Proceedings of the IEEE Military Communications Conference*. IEEE, 1440–1446.

[19] Grace Alexandra Lewis. 2016. *Software Architecture Strategies for Cyber-foraging Systems*. Ph.D Dissertation, Carnegie Mellon University, Pittsburgh, PA.

[20] Jia-Liang Lu and Fabrice Valois. 2006. Performance evaluation of 802.11 WLAN in a real indoor environment. In *Proceedings of the IEEE International Conference on Wireless and Mobile Computing, Networking, and Communications (WiMob'06)*. IEEE, 140–147.

[21] Longjie Ma, Jigang Wu, and Long Chen. 2017. DOTA: Delay bounded optimal cloudlet deployment and user association in WMANs. In *Proceedings of the 17th IEEE/ACM International Symposium on Cluster, Cloud, and Grid Computing*. IEEE Press, 196–203.

[22] Yuyi Mao, Jun Zhang, S. H. Song, and Khaled Ben Letaief. 2016. Power-delay tradeoff in multi-user mobile-edge computing systems. In *Proceedings of the IEEE Global Communications Conference (GLOBECOM'16)*. IEEE, 1–6.

[23] Antti P. Miettinen and Jukka K. Nurminen. 2010. Energy efficiency of mobile clients in cloud computing. *HotCloud* 10 (2010), 4–4.

[24] Anwesha Mukherjee, Debashis De, and Deepsubhra Guha Roy. 2016. A power and latency aware cloudlet selection strategy for multi-cloudlet environment. *IEEE Trans. Cloud Comput.* 99 (2016), 1–14.

[25] Yucen Nan, Wei Li, Wei Bao, Flavia C. Delicato, Paulo F. Pires, Yong Dou, and Albert Y. Zomaya. 2017. Adaptive energy-aware computation offloading for cloud of things systems. *IEEE Access* 5 (2017), 23,947–23,957.

[26] Yucen Nan, Wei Li, Wei Bao, Flavia C. Delicato, Paulo F. Pires, and Albert Y. Zomaya. 2018. A dynamic tradeoff data processing framework for delay-sensitive applications in cloud of things systems. *J. Parallel and Distrib. Comput.* 112 (2018), 53–66.

[27] Deepsubhra Guha Roy, Debashis De, Anwesha Mukherjee, and Rajkumar Buyya. 2016. Application-aware cloudlet selection for computation offloading in multi-cloudlet environment. *J. Supercomput.* (2016), 1–19.

[28] Mark Ryan. 2005. *Calculus Workbook for Dummies*. Wiley Publishing, Inc.

[29] Swetank Kumar Saha, Pratham Malik, Selvaganesh Dharmeswaran, and Dimitrios Koutsonikolas. 2016. Revisiting 802.11 power consumption modeling in smartphones. In *Proceedings of the 17th IEEE International Symposium on World of Wireless, Mobile, and Multimedia Networks (WoWMoM'16)*. IEEE, 1–10.

[30] Mahadev Satyanarayanan, Grace Lewis, Edwin Morris, Soumya Simanta, Jeff Boleng, and Kiryong Ha. 2013. The role of cloudlets in hostile environments. *IEEE Pervas. Comput.* 12, 4 (2013), 40–49.

[31] Jiafu Tang, Chongjun Yan, Xiaoqing Wang, and Chengkuan Zeng. 2014. Using Lagrangian relaxation decomposition with heuristic to integrate the decisions of cell formation and parts scheduling considering intercell moves. *IEEE Trans. Automat. Sci. Eng.* 11, 4 (2014), 1,110–1,121.

[32] Song Wu, Chao Niu, Jia Rao, Hai Jin, and Xiaohai Dai. 2017. Container-based cloud platform for mobile computation offloading. In *Proceedings of the IEEE International Parallel and Distributed Processing Symposium (IPDPS'17)*. IEEE, 123–132.

[33] Zichuan Xu, Weifa Liang, Wenzheng Xu, Mike Jia, and Song Guo. 2016. Efficient algorithms for capacitated cloudlet placements. *IEEE Trans. Parallel Distrib. Systems* 27, 10 (2016), 2,866–2,880.

[34] Hong Yao, Changmin Bai, Muzhou Xiong, Deze Zeng, and Zhangjie Fu. 2017. Heterogeneous cloudlet deployment and user-cloudlet association toward cost-effective fog computing. *Concurr. Comput.: Practice and Exper.* 29, 16 (2017).

[35] Chongyu Zhou, Chen-Khong Tham, and Mehul Motani. 2017. Online auction for truthful stochastic offloading in mobile cloud computing. In *Proceedings of the IEEE Global Communications Conference (GLOBECOM'17)*. IEEE, 1–6.

[36] Qiliang Zhu, Baojiang Si, Feifan Yang, and You Ma. 2017. Task offloading decision in fog computing system. *China Comm.* 14, 11 (2017), 59–68.

**Author Queries**

Q1:   AU: Please supply the CCS Concepts 2012 codes per the ACM style indicated on the ACM website. Please include the CCS Concepts XML coding as well.

Q2:   AU: In Table 1, "the category belong the tasks" written as meant?

Q3:   AU: Please specify figure number instead of "the last figure."

Q4:   AU: Please clarify: "more the remote processing capacity is important compared to the local processing capacity more the decision is to offload task."

Q5:   AU: Please specify figure number instead of "the last figure."

Q6:   AU: Please supply publication source for Reference 11.