

Context-Aware Multi-User Offloading in Mobile Edge Computing: A Federated Learning-based Approach

Ali Shahidinejad^{1*}, Fariba Farahbakhsh¹, Mostafa Ghobaei-Arani¹, Mazhar Hussain Malik², Toni Anwar³

¹Department of Computer Engineering, Qom Branch, Islamic Azad University, Qom, Iran

*E-mail: a.shahidinejad@qom-iau.ac.ir

Abstract:

Mobile edge computing (MEC) provides an effective solution to help the Internet of Things (IoT) devices with delay-sensitive and computation-intensive tasks by offering computing capabilities in the proximity of mobile device users. Most of the existing studies ignore context information of the application, requests, sensors, resources, and network. However, in practice, context information has a significant impact on the offloading decisions. In this paper, we consider context-aware offloading in MEC with multi-user. The contexts are collected using autonomous management as the MAPE loop in all offloading processes. Also, federated learning (FL)-based offloading is presented. Our learning method in mobile devices (MDs) is deep reinforcement learning (DRL). FL helps us to use distributed capabilities of MEC with updated weights between MDs and EDs. The simulation results indicate our method is superior to local computing, offload, and FL without considering context-aware algorithms in terms of energy consumption, execution cost, network usage, delay, and fairness.

Keyword: Mobile edge computing. Computation offloading. Context-aware. Federated learning

1. Introduction

In recent years, data production in various scientific and industrial fields and the limitation of resources to process has resulted in the need for a rich processing environment outside of the user's equipment. This led to creating the cloud computing environment with almost endless physical and virtual processing resources [1]. In computation offloading to the cloud, we face problems due to the large distance of end-users to the cloud; for example, in cases where we need a real-time response, such as healthcare, waiting for receiving a response from the cloud can pose serious problems. A computational offloading was proposed to solve this problem in the fog environment to create a computational level closer to the end-users [2, 3]. For several years, a new trend was emerging and putting cloud computing functions on the network's edge. One of the incentives for this approach is the mass production of network EDs (including Wi-Fi router and access point station). Due to the significant processing power of these devices, high-throughput and delay-sensitive functions can be implemented. This process model is called mobile edge computing (MEC) [4]. This technology is developed by the European telecommunications standards institute (ETSI) [5]. The main focus of the MEC is on radio access networks (RANs) in 4G and 5G cellular networks.

MEC has advanced features like latency, user proximity, high bandwidth, and location awareness [6]. These enable MEC to run many new types of applications and multi-region services, such as business and health, augmented reality, video streaming services, and more [7]. At MEC, the user's distance is much

closer than the user's distance to the cloud. One of the key technologies of MEC is computational offloading. This can be examined from both single-user and multi-user aspects. In a single-user computational state, at any given moment, a user can offload the computing task. In contrast, in multi-user computational offloading, multiple users are allowed to move their tasks to other computing layers simultaneously. As a difference between these two states in the multi-user offloading, one module with different data related to different users can be offload to EDs or the cloud. Since in the single-user mode, each module has got specific data. **Therefore, it can be offload by the user. In the single-mode, there is no need for any data management for users.** Although many works have been done in computational offloading in recent years, the concept of context-awareness has been used very limitedly [8] in past research. **Our meaning about context is using the properties of offloading, application, mobile, sensors, network and media, and resources.** The context in computing offloading decisions will be very influential because of mobile conditions like location, network status, and available computing resources [9].

One of the issues that arise in offloading is intelligent tools to detect current or underlying conditions and implement context-based behavior [10, 11]. This ability can be referred to as context-awareness. As soon as they make a network available, they perform the offload without considering whether the offloading is in their favor or not. The computational burden is not always beneficial to obtain the required level of efficiency and benefit offloading. **Here we are investigating to improve the delay and energy consumption by the proposed offloading method.** The distributed nature of MEC requires an appropriate offloading method. For this purpose, the FL can be useful in this regard. FL can coordinate the training process among multiple MDs. The DRL technique can solve the offloading problem. **The DRL technique is very efficient in finding the optimal offloading policy in MEC. Since DRL needs much processing,** thus the DRL agent has to be carefully designed and implemented. Some challenges [12] and their solutions in FL are as follows.

- The whole training dataset is not accessible. This challenge is created for the nature of distributed computation, and it can provide the privacy of data for all users.
- Slow and unstable communication. Based on the proposed approach, MDs are not completely dependent on EDs. As some nodes become offline, only the weights are less trained or updated later, but the task performing or offloading is done continuously.
- The trade-off between privacy guarantees and system performance. The computation tasks can be encrypted by a fast and trusted cryptography algorithm in IoT.
- Interference among MDs (The MDs may be geographically close to each other. This introduces an interference issue when they update local models to the server. As such, channel allocation policies may need to be combined with the resource allocation approaches to address the interference issue). DRL can be considered to model the dynamic environment of MEC and make optimized decisions.
- Comparisons with other distributed learning methods. Some methods use neural networks up to a cut layer, or others ignore to transmit weights to an aggregating server. FL has a more straightforward implementation since the participants and the FL server run the same global model in each cluster.
- Learning convergence. We improve this challenge with the loss function, as mentioned in the DRL algorithm.

- Size of model updates. The combinations of weights and contexts help us to reduce the size of model updates.

Because real systems and environments are multi-user and not single-user, we use multipurpose computing offloading in our approach. Since these users are located in different locations and conditions, thus the offloading decision should be made with the knowledge of context and the existing conditions. This strategy solves the problems shown above, and to improve the service efficiency in computational offloading, we propose a multi-user conditioned MEC system that changes the conditions of a mobile computing resource. In this research, after presenting a three-tier architecture (IoT devices, edge servers, and cloud), the desired content is collected using autonomous management (MAPE control loop) defined at the edge level (Monitor phase), where we consider some important contexts in this area, using application context, mobile devices, sensors, networks, edge servers, and media. These contexts are analyzed (Analysis phase) and to help make decisions about offloading. These contexts send to our context-aware algorithm then a subsystem (Planning phase) executes the offloading instruction (Execution phase). The question that can be asked here is whether to use the concept of context-awareness in computational offloading in a multi-user MEC. As the context information is exchanged, the FL is implemented, and the updated weights related to the DRL algorithm are shifted between MDs and EDs. Our key contributions in this paper are as follows:

1. We provide a MAPE control loop on the MEC architecture to decide whether to run local or offload computations to edge or cloud. This loop executes in the lifetime of the network and updates all parameters in the problem space. Also, we use context information of the application, sensors, resources, edge servers, and network. These updated contexts improve the offloading process.
2. For optimal use of the distributed capability of MEC, we present an FL-based offloading algorithm. It uses the DRL to train the MDs and sends the updated weights to EDs and the cloud. It causes lower data transmission from MDs to EDs and protects the users' information.
3. The proposed approach evaluates based on some metrics: energy consumption, execution cost, network usage, delay, and fairness. The results show our proposed method outperforms the original, offloading, and FL without context methods.

The rest of this paper is organized as follows. In section 2, related works are summarized. The system model and network architecture are presented in section 3. In section. 4, we explain our offloading algorithms in detail. In Section 5, the evaluation results of our proposed algorithms are presented and compare with other methods. Finally, in Section 6, the conclusion is discussed, and suggestions are made for future work.

2. Related works

In recent years, many studies have been performed about MCC [13, 14] and MEC [4]. We classify these works to multi-user and context-aware offloading as follows. Also, we collect and analyze some researches about FL.

2.1. Multi-user offloading

[There are some research works about multi-user offloading](#) [15,16,17,18]. Here, we mention these papers

based on their objectives and methods. Researchers studied different objectives such as energy consumption [19,44], computation delay [20], QoS (Quality of service), latency, and accuracy [21,22]. According to [16], as with cloud services such as PaaS (platform as a service), IaaS (infrastructure as a service), SaaS (software as a service), cloud computing offloading is also considered as a service (OaaS) in cloud computing. Unlike the client-server method, which the client always requests from the server for the result of a computational task, in the computational offloading method, only when it is needed. The proposed method captures and records user preferences, the current status of devices such as battery level, network bandwidth, CPU speed, free memory, and so on. The simulation results show that computation offloading to a more robust device can improve runtime instead of executing close to the user device.

Some researchers minimized the cost under constraints and solve the offloading problem in multi-user MEC by backtracking, genetic algorithm, and greedy strategies [18]. Paper [23] investigates the computational offloading with an efficient energy scheme in a multi-user fog computing system. In this paper, queuing is used to model the execution processes on mobile and fog devices. The problem of efficient energy optimization is formulated to minimize energy consumption conditional on delay constraints. A distributed algorithm called ADMM (based on the periodic multiplier method) is presented to solve the formulated problem. The simulation shows higher performance than other existing designs. The authors in [44] solved a multi-objective scheduling problem to optimize time and energy consumption. They could improve the objectives by a whale optimization algorithm in the MCC. Paper [45] also worked on the energy consumption and also cost for computation offloading of workflow applications in MEC. This research has been presented by a Non-dominated Sorting Genetic Algorithm (NSGA). The results were better than no offloading and cloud offloading methods.

It has been argued in [24] that although computational offloading can reduce power consumption on mobile devices, it may delay further execution, including sending time between mobile devices and cloud servers. According to theoretical analysis, a multi-objective optimization problem is formulated with reducing energy consumption, execution delay, and payment cost, by finding the optimal computational offloading and transmission power for each mobile device.

The results show decreasing in the mentioned objectives.

In [25], a mixed-integer linear programming (MILP) optimization model was used. This paper considers two types of cloud patches: the local cloud patch and the global cloud patch, which have higher capabilities. The model presented in this paper reduces energy consumption while imposing a significant amount of delay.

Researchers in [26] provides some disadvantages of cloud processing such as high latency and unstable QoS (Data dissemination, routing between mobile devices, and cloud servers). Assuming different real-time computing tasks on different devices, each task is decided to either run locally on the device itself or be offloaded to one of the edge servers or the cloud server. This paper examines low-complexity computing offloading policies to minimize the quality of MEC network service assurance of mobile devices' power consumption. Their method is superior to other compared approaches.

2.2. Context-aware offloading

Considering context information in the offloading problem is done based on different objectives and network architectures such as energy-saving and execution time [27], and latency [28] in MCC [29,30] and MEC. The paper [31] proposed a framework that supports mobile applications with computational offload capability for aware conditions. First, a design pattern was proposed to enable the application to be offloaded on demand. Second, an estimation model was presented to select the appropriate cloud source for offloading automatically. The third is a framework implemented on both the server and client-side. It includes three modules: service selection module, computational offload, and runtime management. The evaluation results and comparison with traditional offshore samples show that the proposed approach can improve runtime and power consumption for highly computational applications.

In [32], an offloading middleware is presented to the aggregate cloud by considering energy level, processing power, runtime, and network bandwidth. In this paper, the resource allocation problem is formulated as a multi-objective optimization that aims to optimize the completion time of the task and the energy consumption of all participating mobile devices by satisfying the task boundary. An NSGA-II is used to obtain the beam solution set. Second, a multi-attribute decision-making (MADM) technique is used to determine the best compromise solution based on the entropy technique and weighting for a priority order. Evaluation Results show that the proposed method manages well the compromise between completion time and energy consumption.

The researchers in [46] analyzed the context-aware energy optimization for services on MDs. Their evaluation was based on three supervised machine learning methods as naïve Bayesian, decision tree, and random forest. They provided this result that using the machine learning method is better than others for reducing the service execution time and the energy consumption in MCC.

In [33], a fault-tolerant aware mobility offloading (MAFO) approach is presented that collects network information and user mobility over time and uses the Markov chain of the user's visited networks in different possible paths. It also predicts the stoppage time of each network based on user mobility. The evaluation results show that improvements in time and energy consumption. Authors in [34] have suggested a framework called Thinkair that simplifies developers' work to migrate their smartphone apps to the cloud. It uses the concept of smartphone virtualization in the cloud and provides method level computational offload. It focuses on the resilience and scalability of the cloud and enhances the power of cloud computing by implementing a parallelization approach using multiple virtual machines. The results show that better performance and lower power consumption than similar non-parallel methods. In [35], a framework is considered to decide whether to offload a given method to cloud servers. In this paper, a field-aware decision-making algorithm is designed, implemented, and evaluated called CADA, which uses user contexts and historical metrics to optimize the performance of mobile devices with various optimization criteria such as short response time. The evaluation results demonstrate the high accuracy of CADA algorithm prediction and improving response time and energy consumption.

2.3. FL-based offloading

Using cooperative models have shown good performance in the IoT devices [36]. FL is a cooperative-based

method that can be used in MEC. Here we first try to present a conceptual view of FL and then provide some offloading problems that are solved by FL.

FL allows users of devices to collaboratively train a shared model while keeping data privacy on devices. Thus, FL can be used as an enabling technology for ML model training at MEC. In fact, each device can process its task by a learning model. This can happen on all devices. After that, all devices can share their experience together. As a result, we will have a global model by aggregating all learning models. This is very important that in this cooperation, any private data not transferred between devices.

Generally, there are two main entities in FL as N data owners as $\{1,2,\dots,N\}$ and the model owner (FL server). According to Fig. 1, in the initialization as step 1, the FL server specifies the hyperparameters of the global model and the training process, e.g., learning rate. Then, in step 2, each data owner i (Mobile device) train a local model w_i and send it to the FL server (Controller). In step 3, all collected local models in the FL server are aggregated $w = \bigcup_{i \in N} w_i$. Steps 2 and 3 are repeated until the global loss function converges or a desirable training accuracy is achieved [12].

Since the offloading methods in MEC need a real distributed algorithm; thus, FL is an excellent way to this purpose. In [37], the authors presented an FL-based offloading in MEC. DRL algorithms executed in MDs and updated weights transfer between MDs and ED. Their used parameters were energy consumption and transmission time. The results show a better result than centralized DRL. Because in centralized DRL, the tasks are waiting in a queue to get resources of devices. It might a number of tasks be dropped due to insufficient resources. Nevertheless, FL can offload some tasks to other devices. This causes a lower drop task in devices.

A group of researchers proposed an aggregation model of EDs in the Cloud by FL. They used the difference of convex functions (DC) representation for sparse and low-rank function [38]. It is demonstrated that the novel method was able to select more devices than other benchmark approaches. The paper [39] is presented based on a distributed DRL in MEC for caching and communication operations. This research includes three parts, information collecting, cognitive computing, and request handling. The results show some improvements in utility for the user equipment. This study is well-done; however, it could have evaluated and compared with other state-of-the-art methods.

FL is also used to make decisions about computation offloading and energy allocation in MEC [40]. Here, a DRL-based algorithm is proposed to maximize the expected long-term utility. This method has better results than the centralized and greedy-based offloading algorithms.

As stated in the introduction section's contribution list, the main idea of this research is to use MAPE and FL-based offloading algorithm. Context information has been used in the previous works. We have tried to offload the modules with these contexts in the controlled MAPE loop with our distributed algorithm to improve the mentioned objectives. We try to compare our results with new researches. Also, federated learning is very close to the distributed learning paradigm. In previous works, DRL or DL has been used in each device, and devices have not any relation together. We solve this problem with federated learning.

The summary of offloading algorithms is categorized by technique, objectives, architecture, pros, and cons in Table 1.

Table 1: Summary of offloading algorithms (ET: execution time. Arc: Architecture).

| Ref | Technique | Objectives | Arc. | Pros | Cons |
|------|---|--------------------------------|------------|--|---|
| [31] | FL-DRL | Energy and transmission time | MEC | Cooperative model with improve bandwidth crowding | Using only one ED |
| [32] | FL-DC | Accuracy and number of devices | Edge-Cloud | Selecting maximum devices | Ignoring edge computing objectives |
| [16] | Backtracking, genetic the algorithm, and greedy | Cost | MEC | Minimizing cost | Time-consuming method |
| [34] | FL-DRL | long-term utility | MEC | Analysis DRL parameters and energy efficiency | Ignoring privacy evaluation |
| [23] | LP, deep learning | ET, Energy | MEC | Improve network service quality and reduce mobile device power consumption | Long run time of the proposed algorithms |
| [21] | Queuing | Delay, cost, and energy | Fog | Increasing performance with testing transmission time in the worst condition | Non-optimal method in global goals |
| [14] | Oaas and matching algorithm | Cost | Fog | Real tools for offloading | Ignoring computational offloading time delay |
| [20] | Queuing and ADMM | Delay, energy | Fog | Cell-level alignment, offloading with minimized power consumption | Solving the problem only up to the level of servers and override local implementation and dump computing load to the cloud |
| [22] | MILP | Energy | MEC | MILP for force optimization of large-scale and MEC without delay limits defect | The routing process of requests through the hierarchical cloud patch network only tolerates a significant amount of delay, limited capabilities of the cloud patch system |
| [26] | Game theory, context (Network connection, runtime status) | Energy, ET | MCC | Decreasing ET and energy consumption | High complexity and no comprehensive method |
| [28] | Markov chain, contexts (User mobility, device, program, cloud server, mobile network) | Energy, ET | MCC | Fault tolerance, low energy consumption, and ET | Non-comprehensive evaluation |
| [27] | OMMC, NSGAI, TOPSIS, contexts (Device context, processing power, | Delay, energy, deadline | MEC | Good result with the trade-off between completion time and energy consumption | Ignoring method's complexity |

| Ref | Technique | Objectives | Arc. | Pros | Cons |
|------|--|------------------|------|---|----------------------------------|
| | | | | | CPU usage, network bandwidth) |
| [29] | Thinkair and contexts (Hardware, software, network condition, energy estimation model) | Cost, ET, energy | MCC | Resource allocation based on requests, parallel reconstruction of VMs | Unsuitable for IoT applications |
| [30] | CADA and contexts (Signal strength, transmission time, time-of-day, and location) | Energy, ET | MEC | Using daily time for offloading, decreasing energy consumption and ET | Weak energy model |

3. System model and problem formulation

In this section, we present the architecture and system model. Fig. 1 shows the three-layer architecture of our proposed system. This model includes some sections as follows.

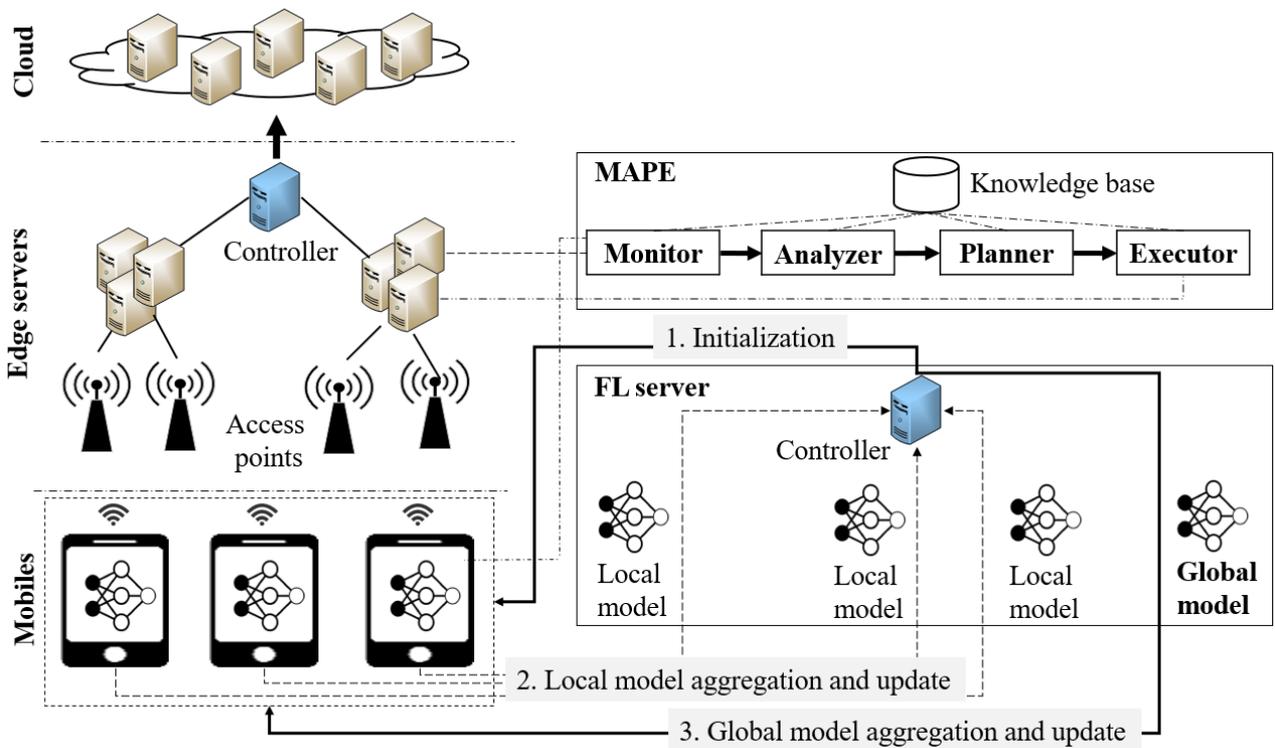


Fig. 1: The architecture and system model

3.1. IoT

The IoT component is at the very bottom of the architecture, including communication devices that are connected through heterogeneous networks. In general, it aims to collect and process data through IoT devices to extract patterns and discover patterns or to perform predictive analysis or optimization and make smarter, more timely decisions. However, mobile devices are considered because of the problem of offloading. The data is first collected by the IoT devices described in this framework, and each user sends

their requests to the queue according to the data collected.

3.2. Controller

This component plays the role of the master node in our hierarchical model. This component is located at the edge command center and is at the top edge of the edge. This component itself is a robust edge resource that manages resources and sourcing for requests from the lower layer. That decides whether the request will be executed on the same edge layer or delivered to the cloud layer based on the users' context and existing resource conditions (existing edge server features) or if the request is to be in the edge layer. Run this component to make it a proper edge server component that can execute the request if the submission is sent and must be moved to the cloud, the cloud component transmits the request to the gateway.

3.3. Edge server

The edge layer component consists of several edge server components that play the role of a slave in our hierarchical approach. These components send information about their processing and storage capabilities to the edge server controller component. This component selects one of them to execute the desired request by matching the context information and resource capability.

3.4. MAPE

This component, which is the main component of our framework, includes components (monitor, analyze, plan, and execute). This component collects available conditions and resources available and extensions of the IoT devices. Mapping these tasks examines the available resources and decides whether to execute them on the edge layer or offload the super layer's computation. This component rests on the edge server controller component. Our context-aware algorithm is implemented in this component. To achieve autonomous computing, IBM has proposed a reference model for autonomous loops, known as the MAPE-K loop, which has four components (monitor, analyze, plan, execute) [41]. These four components, under common knowledge in the MAPE-K autonomous loop, an intelligent agent understands its environment through sensors and uses these perceptions to determine what actions to take in the environment. All four phases, covered by common knowledge they connect and exchange information.

3.4.1 Monitor

In the monitoring phase, the properties of the environment are recorded by the sensors. The data is first received through sensors and intelligent equipment, and according to the data received, a request for execution is made. Our systems and priorities are categorized and, if they are in accordance with the circumstances, these requests are planned for implementation in the third phase (planning).

3.4.2 Analysis

The analysis phase deals with the processing of metrics collected from the monitoring phase, and by processing these metrics, it obtains data on the status of the current productivity of the system and forecasts of future needs so that, if necessary, an appropriate response is obtained. In the analysis phase, warnings and threats are considered. Any violation of the level of needs defined in the analysis phase is considered.

3.4.3 Planning

In the third phase, as planning, based on the tables created in the previous two phases, an appropriate decision is made that leads to offloading or local computing.

3.4.4 Execution

The fourth phase, as execution-only, executes the planned third phase (execution) instructions. In fact, it is responsible for executing the programs approved by the analysis phase. We propose a hierarchical model for the proposed system, in which the edge layer plays a node (Master) in which all four phases of the MAPE loop are implemented, and the other nodes play a role (Slave). With this smart, autonomous solution, decentralized collections are managed in a centralized system. Integrating a centralized and distributed strategy can be important as an innovative strategy. Autonomous loop computing (MAPE) decision-making autonomously leads to better management of resources, reduced response time to heavily time-dependent applications and requests, and reduced system latency.

3.5. Cloud

When requests from the edge server controller component are decided to go to the cloud layer, they are sent directly to the cloud gateway and through it go into the cloud and, depending on which server the requests are to be processed, go to the server's dedicated queue. Finally, when it is time to move to the servers to do their job, depending on request processing and storage type.

3.6. Case study

The VR-GAME (Virtual reality Game) application is a human-based game. According to the workflow of this application, EEG signals send to the client module. The client module sends consistent data to the concentration calculator module. The client module updates the game display to the player. The coordinator module gathers and distributes measured concentration among players [42].

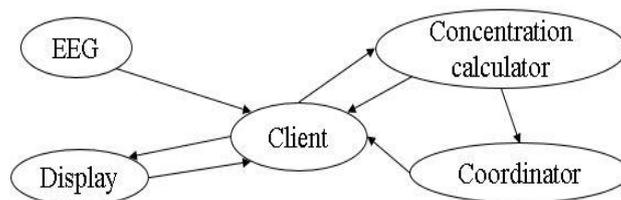


Fig. 2: Application of virtual reality game

The EEG value could be used to determine the interval between two sensed signals. Based on the application part of Figure 2, the EEG sensor, display actuator, and client module are placed in the mobile device. The concentration calculator and the coordinator modules can be placed in the EDs or cloud.

The main problem of this paper is the offloading of modules as $\{M_1, M_2, \dots, M_k\}$ to edge servers as $\{ED_1, ED_2, \dots, ED_n\}$ or Cloud. The problem formulation is explained as follows. The symbols used in this paper are defined in Table 2.

Table 2: Symbol definition

| Symbol | Definition |
|--------------|---|
| μ | CPU cycles for processing the task |
| g_i | Number of energy units |
| ω | Commonly adopted effective switched capacitance |
| s_i^e | Channel gain between the MD and an ED in epoch i |
| f_i^{tr} | Transmission power |
| d | Transmission data size |
| P_i | Power of i th ED |
| U_i | Utilization of i th FD |
| U_i^p | Utilization in the previous updates time |
| T_1 | The time frame of datacenter |
| T_2 | The time frame of host |
| TD | Difference between current and last process time |
| U_{MIPS} | Utilization of MIPS |
| TAM | Total allocated MIPS |
| E_c | Current energy consumption |
| T_n | Current time |
| P_h | Host power in last utilization |
| N | Number of FDs |
| C | Execution cost |
| S_c | System clock |
| T_{lu} | Last utilization update time |
| R_M | Rate per MIPS |
| ω | Number of processors in a host |
| U_l | Last utilization |
| $MIPS_k^A$ | Allocated MIPS of k th processor in the host |
| m_1 | Number of all processors and allocated processors |
| m_2 | Number of allocated processors |
| T_{ED} | ED's execution time |
| $MIPS_k$ | MIPS of k th processor in the host |
| TMA | Total allocated MIPS of the host |
| L_i | Total latency |
| S_i | The total size of i th tuple |
| N' | Total number of tuples |
| T_{max} | Maximum simulation time |
| T_{st} | Tuple start time |
| T_a | Average CPU time of the tuple type |
| ET | Emitting time of a tuple |
| T_s | Sending time of a module to another module |
| N_{RT} | Number of receipt tuple types |
| Q | The number of devices contributed to the offloading |
| x_j | Energy consumption of j th device |
| A | Selected action by the agent |
| R | Reward value |
| α | Learning rate |
| γ | Importance of the next rewards |
| $Q(S,a)_m^i$ | Q update value |

3.7. Local execution time

If the required resources in MDs are provided, thus we can calculate the time consumption locally. This value calculates in each episode i . In Eq. (1), μ is CPU cycles for processing the task. Finally, ω is the commonly adopted effective switched capacitance that depends on the architecture of chips [37,47]. According to [47], ω can be given by $\sum_{i=1}^N \alpha_i * C_{L_i} * \Delta V_i$, where C_{L_i} is the physical capacitance, α_i is the activity weighting factor, each averaged over the N nodes in the circuit. Also, ΔV_i is the voltage change. It is required to explain, the ω value is calculated by the simulator.

$$T_i = \frac{\mu}{\sqrt{\frac{1}{\omega * \mu}}} \quad (1)$$

3.8. Data rate between MDs and EDs

If the MD wants to communicate with an ED, and a wireless link is established for them. The achievable data rate calculates by Eq. (2). Here, A is the power of interference plus noise. S_{i^e} is the channel gain between the MD and an ED in epoch i [37]. **This channel gain is static and independently taken from the state of MD.**

$$DataRate_i = \omega * \log_2 \left(1 + \frac{S_{i^e} * f_i^{tr}}{A} \right) \quad (2)$$

The transmission power is calculated by f_i^{tr} as Eq. (3). BW_i is the bandwidth of ED in epoch i. where d is the transmission data size required for offloading a module.

$$f_i^{tr} = \frac{BW_i}{d} \quad (3)$$

3.9. Edge server's power consumption

The power consumption of each ED is presented here. According to this equation, an edge server with more power than the rest of the edge servers is a candidate for an offloading destination.

$$P_i = P_i^c + T1 + T2 \quad (4)$$

In Eq. (4), P_i^c is the current power of FD and T1 is the energy consumption of the datacenter in the current time.

$$T2 = U_i^p + \frac{U_i - U_i^p}{2} * TD \quad (5)$$

In Eq. (5), U_i is the utilization of ith FD. U_i^p is the utilization in the previous updates time. TD is the time difference between the current time and the last process time.

$$U_i = \frac{U_{MIPS}}{TAM} \quad (6)$$

The total allocated MIPS of all processing elements is updated as EQ. (7).

$$TAM = \sum_{i=1}^N \sum_{j=1}^M PEM_{ij} \quad (7)$$

Where TAM is the total allocated MIPS of an ED that is less than or equal to MIPS of that ED ($TAM \leq EDMIPS$).

$$U_i^p = \frac{U_{MIPS}^p}{TAM} \quad (8)$$

In Eqs. (6) and (8), U_i^p is the utilization value in the previous time, U_{MIPS} is the utilization of MIPS, and TAM is the total allocated MIPS [42].

3.10. Edge server's execution time

MIPS calculates the runtime of modules in edge servers. The number of commands that any edge server can

handle given its current workload is considered its current capacity. Therefore, according to Eq. (9), each module can capture and run it at T_{ED} , where MIPS is the million executable operations that an edge server can run per second.

$$T_{ED} = \frac{1}{MIPS} \quad (9)$$

3.11. Edge server's bandwidth

Each ED includes some hosts as follows.

$$\{Host_1, Host_2, \dots, Host_n\} \in ED_i \quad (10)$$

The main properties of hosts are RAM, bandwidth, storage, and PEs. In a host, as Eq. (11):

$$BW_{Lower} \leq \sum_{i=1}^N BW_i \leq BW_{Upper} \quad (11)$$

Where BW_{Lower} is the lower bandwidth, and BW_{Upper} is the upper bandwidth of each ED. BW_i is the bandwidth of i th host. N is the number of hosts. The total bandwidth of all hosts in each ED is between BW_{Lower} and BW_{Upper} .

4. The proposed approach

As stated above, our goal is to apply the concept of context knowledge to a multi-user mobile edge computing system. The proposed framework for the context-aware system can be described as follows. In this system, we have two types of variables: independent variables and dependent variables. Independent variables are all input variables that the system receives in the form of transactions and does not interfere with their calculation, such as the types of fields that surround the environment. Dependent variables are variables that are obtained by using independent variables as inputs to the proposed system. Delay and energy consumption are those variables. The following sections describe the various tasks in MAPE.

4.1. MAPE

The MAPE control loop consists of four phases the monitor, analyze, plan, and execute. We explain them as following steps:

4.1.1 Monitoring: This section monitors and collects input modules. This is basically the context monitor component of our system. The inputs include all requests received from IoT devices and fields collected from the environment which enter modules. The request is received with a unique identifier. This request can be either computation or data. In this section, independent parameters such as QoS and SLA are also monitored and written to the knowledge database. We consider the contexts to include application, mobile device, sensors, network, and media.

- **Offloading contexts:** Request id, requester name, sensitivity type (resource-based or time-based), QoS, and SLA requirements. Based on context information, the QoS depends on data rate between MDs and EDs as Eq. 2, edge power consumption an Eq. 4, and edge server's execution time introduced in Sec. 3.6.4.

- **Application contexts:** Total executed modules, runtime, allocated memory, priorities of modules, and source type.

$$Priority_i \leq Priority_j \quad (12)$$

According to Eq. (12), the i th module has got more priority than the j th module if the module is before the j th module in the application's workflow.

- **Mobile contexts:** Average frequency of CPU, average CPU usage, and battery level. The CPU usage depends on the MIPS of each CPU.
- **Sensor contexts:** Sensor id, location, latency, destination module, tuple type, and transmission time.
- **Network and media context:** Cellular communication and bandwidth mode, Wi-Fi communication mode, cellular connectivity signal, and Wi-Fi.
- **Resource contexts:** Resource state, identification, memory, and storage.

All of these fields are stored in our knowledge bank's context database to be used later in the computing offloading operations.

4.1.2 Analysis: This component deals with the processing of metrics collected from the monitor component, and by processing these metrics, it obtains data on the status of current system productivity and forecasts of future needs. In this phase, QoS and SLA are considered. If a resource is assigned to a computing request, and this results in a breach of service quality, the analyst must detect this and issue the necessary alert. The second phase of the loop performs such analyzes. This department has a close relationship with the knowledge bank and is constantly exchanging information with it. The analysis phase's output contexts are resource id, offloading request-id, QoS, and SLA types.

4.1.3 Planning: This part contains the decision module of our system. Using the information from the previous section, this section makes the final decision on whether to offload or execute locally through the knowledge bank. This component decides whether to offload and if so, how to send the task to the appropriate infrastructure under the current circumstances. The decision module includes two components of the cost estimation module and the context-aware decision algorithm. This method finds the best destination for offloading the requester modules to edge-server or cloud. We present two offloading methods as MUCAO and FLUCO as following.

4.1.3.1 MUCAO

Our first method for finding the best destination for offloading is based on a conditional technique. Algorithm 1 includes two sections as initialization and MAPE. In the initialization section, firstly, mobile devices, cloud, applications are created. Secondly, edge servers are built based on the number of departments and mobiles. Finally, the application is submitted to the edge server controller, and iFogsim is started. In the MAPE section, four phases execute continuously. In the monitoring phase, the context of modules, sensors, tuples, network interfaces, mobile devices, cloud, and edge servers are collected. In the analysis phase, the cost of execution in the local device, edge server, and cloud is calculated, and the network interface state is checked. In the planning phase, the availability of local devices, edge servers, and cloud are checked. Finally, the current module offloads to the best destination in the execution phase, which is obtained in the planning phase or locally executed.

Algorithm 1 MUCAO

```
1: Create Mobile devices, Cloud, application (Modules, Edges, Tuples, Workow).
2: for i = 1 to DepartmentMax do
3:   for j = 1 to MobileMax do
4:     Edge server (Node name, MIPS, Ram, Storage, upper BW, lower BW, busy
power, and idle power).
5:   end for
6: end for
7: Submit applications.
8: Start iFogsim.
9: while Modules enter from MDs do
10: Monitor:
11:   Collect context of modules, sensors, tuples, network interfaces, mobile
devices, cloud, and edge servers.
12: Analyze:
13:   Calculate cost of execution in local device, edge server, and Cloud.
14:   Check network interface state.
15: Plan:
16: if Available(Local device) then
17:   if Available(Cloud) then
18:     Decision = MinCost(Local device,Cloud).
19:     Break.
20:   else
21:     Decision = local device.
22:   end if
23: else
24:   if Available(Wi-Fi) then
25:     if Available(Edge server) & Available(Cloud) then
26:       Decision = MinCost(Local,Edge,CLoud).
27:     end if
28:   end if
29: end if
30: Execute:
31: Offload module based on Decision.
32: end while
33: Stop iFogsim.
```

4.1.3.2 FLUCO

The second proposed approach is based on the DRL. The DRL approach aims to learn the optimal MEC offloading policy from past experience. We try to extend this method to the distributed system. Our offloading algorithm implements in the MDs. The EDs and cloud devices analyze the updated weights from MDs. Then, each MD can decide to offload tasks to the best devices for execution. Here, we define the module offloading by DRL's agent as a tuple:

$$Agent = (M, S, A, Q) \quad (13)$$

In Eq. (13), M is the set of modules' attributes for allocation by agent, S is the set of all environment states, A is the set of actions like local execution, FDs, or Cloud, and Q is the quality function that learning algorithm can select the best destination for module execution by that. These parameters use for the agent's action and calculated by Eq. (14).

$$Q: X^i * A \rightarrow \mathbb{R} \quad (14)$$

Here, x^i is based on Eq. (15), A is the selected action by the agent, and R is the reward.

$$x^i = (z_1(x), z_2(x), \dots, z_n(x)) \in X^z \quad (15)$$

As Eq. (16), each agent can explain modules and environments by Z .

$$Z = M \cup S = \{z_1, z_2, \dots, z_n\} \quad (16)$$

The Q update function is as Eq. (17). $\alpha \in (0, 1)$ is the learning rate, $\gamma \in (0, 1)$ show the importance of the next rewards.

$$Q(S, a)_m^i = Q(S, a)_m^i + \alpha[r + \gamma \max_{a'} Q(S', a')_m^i - Q(S, a)_m^i] \quad (17)$$

According to DRL, the maximum value of $Q(S, a)_m^i$ based on action a is $1 - \epsilon$, and other actions have ϵ probability. Using a greedy policy [43] technique is to avoid local optimum in the learning algorithms. A reward function evaluates agent operations [43]. This function should be able to generate output very quickly so that learning and problem solving can be done without delay. The reward function in the proposed approach is as Eq. (18).

$$Reward = \frac{P_i}{T_i} \quad (18)$$

Where P_i is the available power of i th FD and T_i is the execution time of a module in i th FD. Since power and execution time values are in different ranges, a logarithmic function is used to normalize them in $[0, 1]$.

Thus $P_i = \frac{\log P_i}{10}$ and $T_i = \frac{\log T_i}{10}$.

DRL algorithm:

The pseudo-code of the DRL is as Algorithm 2. The learning algorithm is executed for all modules. Then, the Q-table is initialized by 0. For all episodes, the possible actions and Q values of them are calculated, and the best action is selected by $\arg \max Q$. State S' is transferred to state S . Here, the best destination for each module is selected. After calculation of the reward function, updating the operation of Q and saving episodes in memory is done.

Algorithm 2 DRL

```

1: for m = 1 to LastModuleInQueue do
2:   Initialize Q(S,a).
3:   for i = 1 to EpisodeLast do
4:     Set S to S0.
5:     for j = 1 to SLast do
6:       Select best a by calculation  $\arg \max Q$ .
7:       Action a, visit r and S'.
8:       Calculate  $Q(S, a)_m^i$  as Eq. (17).
9:       Transfer S to S'.
10:    end for
11:   end for
12:   Select a destination device for each module.
13:   Calculate real-time reward.
14:   Save S, r, S', and  $\alpha$  in memory.
15:   Train Q policy by training samples.
16:   Update Q.
17:   Save the current episode in memory.
18: end for

```

Based on our approach, the DRL does not execute in EDs for three reasons.

- 1) The jeopardize of accessing personal data of MDs.
- 2) Encryption algorithms can protect data, but communication with the EDs weaken MDs' privacy.
- 3) Transferring a lot of data from MDs to EDs causes a lot of bandwidth consumption and burden wireless channels. We face another challenge. If the DRL agent runs on each MD, it will consume a lot of energy and time. For solving this problem, we have to explain, our proposed method is not based on the separate learning of each MD, and we use the distribution capability of the MEC. In fact, we propose FL for distributed training DRL agents. As a result, we will save a lot of energy.

FLUCO algorithm:

In FL-based offloading, each ED is a controller to coordinate some MDs. Each MD can execute a DRL agent with less computation burden. FL does three steps:

- 1) Send the DRL agent's parameters from the ED.
- 2) MDs use to download data from EDs for upgrading their model.
- 3) Send updated DRL agent's parameters from MDs to ED (model aggregation). FL works in a parallel manner. This increases the performance of the system. To design an optimal control policy on FL, we have to maximize the expected long-term utility as Eq. (19).

$$G(E, U) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N g(E_i, U(E_i) | E_1 = E) \quad (19)$$

Where E is the network size, U is the system utility, E_i is the initial network size, $g(0)$ is the immediate utility at epoch i that is calculated based on the reward function in DRL. Based on the mentioned approach as Algorithm 2, DRL agents execute in MDs; training is done, local execution or offloading to best ED are done. Finally, trained weights upload to EDs. EDs do not execute DRL agents and only update and aggregate their weights and send them to MDs. The weights aggregate by Eq. (20).

$$W_{r+1} = \sum_{k=1}^{Set_{Last}} \left(\frac{C_r^i}{C_r} * W_{r+1}^i \right) \quad (20)$$

Where W_{r+1} is the weight in the next round, Set_{Last} is the last set of available MDs, C_r^i is the context of i th MD in round r , and W_{r+1}^i is the weight i th module in the next round. This process does continuously. Thus, the computation task executes in MDs or offload to the best ED based on the DRL agent result. One ED is used in some FL-based offloading methods, and model updates of all MDs transfer to ED. This has many challenges for bandwidth, congestion, centralized issues, privacy, etc. We propose multi ED in the MEC and update all EDs according to MDs' download information, which benefits MEC capabilities.

Algorithm 3 FLUCO

- 1: Create Mobile devices, cloud, application (Modules, EDs, Tuples, Workflow).
 - 2: **for** $x = 1$ to $Department_{Max}$ **do**
 - 3: **for** $y = 1$ to $Mobile_{Max}$ **do**
 - 4: Create ED (Node name, MIPS, Ram, Storage, upper BW, lower BW, busy power, and idle power).
 - 5: Initialize the DRL agent with random weights W_0 in the current ED.
 - 6: Initialize the gross training times T_0 .
-

```

7:   end for
8: end for
9: for M = 1 to MDsLast do
10:   Initialize the contexts  $C^{M_0}$ 
11:   Initialize the DRL  $W^{M_0}$ 
12:   Download  $W_0$  from the closest ED.
13:    $W^{M_0} = W_0$ .
14: end for
15: Submit applications.
16: Start iFogsim.
17: while Modules enter from MDs do
18:   for r = 1 to  $r_{Last}$  do
19:     Monitor:
20:     Collect the context of modules, sensors, tuples,
network interfaces, mobile devices, cloud, and edge servers.
21:     Analyze:
22:     Setr = random set of available MDs.
23:     for i = 1 to SetrLast do
24:     Fetch  $W_r$  from ED as  $W_r^i = W_r$ .
25:     Update context  $C_{round}^i$ .
26:     Plan:
27:     Train the DRL agent with  $W_r^i$  on  $C_i$ .
28:     Upload trained  $W_{r+1}^i$  to the closest ED.
29:     Notify the ED the time's  $T_r$  of local training.
30:     end for
31:     for j = 1 to EDLast do
32:     Receive all model updates.
33:     Update  $T_r^j$ .
34:     Aggregate by Eq. (20).
35:     end for
36:   end for
37:   Execute:
38:   Offload modules based on the FDL result.
39: end while
40: Stop iFogsim.

```

The output contexts of the planning phase are Offloading request id, resource type, offloading destination, and considered media for the relationship with a resource.

4.1.4 Execution: The final responsible for executing the commands in the execution section. In this section, computations offload to other machines. This section is closely related to the equipment and resources and stores the latest state of the resources previously mentioned in the knowledge bank for future use. This section includes our task manager component. The task manager collects information such as (method entries, libraries needed to execute the task, the network address of the download location), and puts it into an offloading package. The manager decides to run a hand over the task locally or sends it to the top layers as edge servers or cloud. The output contexts of the planning phase are offloading request-id, resource type, offloading destination, and considered media for relation with the resource.

5. Evaluation

The performance of our proposed methods is presented in this section. The simulation environment in this research is the iFogsim library [42]. This simulator has got classes to implement resource management strategies. We have extended some classes as the ModulePlacementEdgeward for offloading, controller for more output metrics. Also, the VRGAMEFOG class is customized based on the architecture of this paper. We simulate the proposed algorithm and compare the results with other scheduling methods as follows.

- **Original:** In this method, the computations in MDs execute locally. Thus, the computations do not offload to edge servers or the cloud. Since modules continuously need to execute, MDs might not have enough resources. As a result, some modules wait in a queue of resources, and delay will increase.

- **Offload:** In this technique, the destination of tasks or modules is calculated based on the order of edge servers or cloud in the network [15]. Context-aware is not considered in this method. Here, all devices are in a list, and the controller assigns those modules that need resources to the elements of this list in order. This method is not optimal due to ignoring the properties of devices, applications, and network environment. Maybe a device in the first of the list selects for an offload. However, some devices in the middle or last of the list are the best destination for offloading modules.

- **MUCAO:** This algorithm is one of our proposed algorithms. It is an improved offload technique in [15]. We add context-aware to that work. As we presented in the proposed approach section, this method is based on a MAPE loop and uses the execution cost in MDs, EDs, and cloud. Considering context awareness of devices, applications, and network environment leads to find the best device for offloading modules.

- **FLO:** As a state-of-the-art algorithm, FLO is an FL-based algorithm based on DRL [37]. FLO used one ED. Using one ED converts the computing architecture to the cloud. If the number of modules that need resources increase in MDs, just one ED might not answer all offloading requests. As a result, some modules wait for a long time. There are two differences between this work and our proposed algorithm. FLUCO uses many numbers of EDs. Also, we provide the context-aware.

We run the simulation by 3 departments and 4 mobile devices. The comparisons are based on the best results of algorithms with the same configuration for each case study.

5.1. Simulation configuration

Here, we present the VRGAMEFOG application configuration in edges, devices, connection latency, and hosts in Tables 3, 4, and 5, respectively. In Table 3, Pr is the periodicity (mS) of edges.

Table 3: VR game application edge configuration

| Source Module | Destination Module | Pr | Tuple CPU length (B) | Tuple new length (B) |
|--------------------------|--------------------------|-----|----------------------|----------------------|
| EEG | Client | 0 | 3000 | 500 |
| Client | Concentration Calculator | 0 | 3500 | 500 |
| Concentration Calculator | Coordinator | 100 | 1000 | 1000 |
| Concentration Calculator | Client | 0 | 14 | 500 |
| Coordinator | Client | 100 | 28 | 1000 |
| client | Display | 0 | 1000 | 500 |

The host configuration is as follows. The architecture is x86, OS is Linux, Storage is 10^6 B, BW is 10^4 BS, VM model is Xen, the cost is 3 \$, cost per memory is 0.05 \$, cost per storage is 0.01 \$, and time zone is 10. Table 4 shows the parameters of devices including MIPS, RAM (KB), UpBW (Upper bandwidth by kilobyte per second), DownBW (Down bandwidth by kilobyte per second), level in the hierarchical topology, the rate per MIPS, busy, and idle power (Megawatt).

Table 4: Devices configuration

| Device | MIPS | Ram | Uplink BW | Downlink BW | Lv | Rate per MIPS | Busy Power | Idle Power |
|--------|-------|-------|-----------|-------------|----|---------------|------------|------------|
| Cloud | 44800 | 40000 | 100 | 10000 | 0 | 0.01 | 1648 | 1332 |

| | | | | | | | | |
|------------------|------|------|-------|-------|---|---|--------|--------|
| Controller | 2800 | 4000 | 10000 | 10000 | 1 | 0 | 107339 | 834333 |
| ED _s | 2800 | 4000 | 10000 | 10000 | 2 | 0 | 107339 | 834333 |
| MD _{ss} | 500 | 1000 | 10000 | 10000 | 3 | 0 | 8753 | 8244 |

Table 5: Connection latency

| Device Name | Device Name | Latency (mS) |
|----------------------|----------------------|--------------|
| Cloud | Proxy – Server | 100 |
| Proxy – Server | Department (Gateway) | 4 |
| Department (Gateway) | Mobiles | 4 |
| EEG sensor | Mobile | 6 |
| Display | EEG sensor | 1 |

Table 6 shows three different mobile types that have been used in this work. Type A is an Apple iPhone 11, type B is a Samsung Galaxy S10, and type C is a Huawei P30 pro. Their properties include CPU, memory, and battery. In this table, MT is a mobile type.

Table 6: Mobile types

| | Brand | CPU (GHz) | RAM (MB) | Battery (mA) |
|---|--------------------|--------------|-------------|-----------------|
| A | Apple iPhone 11 | 6*2.96 | 4000 | 3110 |
| B | Samsung Galaxy S10 | 8*2.30 | 8000 | 3400 |
| C | Huawei P30 pro | 8*2.03 | 8000 | 4200 |

5.2. Metrics

To analyze our proposed approach and compare it with other offloading algorithms, we consider some metrics like energy consumption, total execution cost, total network usage, delay, and Jain index.

5.2.1 Energy consumption: The energy consumption is calculated by Eq. (21) for all edge servers and cloud when they have serviced the input modules.

$$E = E_c + (T_n - T_{lu}) * P_h \quad (21)$$

We calculate the edge server's energy consumption by the power of all hosts in a certain time frame of execution. Where E_c is energy consumption in the current state, T_n is the current time, T_{lu} is the update time at the last utilization, and P_h is the power of a host in the last utilization. To calculate the total energy consumption, we have to sum all edge servers and the cloud's energy.

5.2.2 Total execution cost:

To obtain the execution cost, we calculate the total MIPS of hosts by the time frame. The time frame is calculated by the difference between the current time of simulation and the last utilization time.

$$Cost = \sum_{i=1}^N \left[C + (SC - T_{lu}) * R_M * U_l * \sum_{k=1}^{\omega} MIPS_k \right] \quad (22)$$

In Eq. (22), N is the number of edge servers, C is the execution cost, SC is the system clock or current time of simulation, T_{lu} is update time at the last utilization, R_M is the rate per MIPS that is different for each inter-module edges, and TM is the total MIPS of the host. U_l is the last utilization (U_i) that is

calculated by Eq. (23). Where $MIPS^A_k$ and $MIPS_k$ are the allocated MIPS and MIPS of the kth processor in the host, and $m1$ and $m2$ are the number of all processors and allocated processors in a host, respectively.

$$U_l = \text{Min}\left(1, \frac{\sum_{k=1}^{m2} MIPS_k^A}{\sum_{k=1}^{m1} MIPS_k}\right) \quad (23)$$

5.2.3 Total network usage

Since tuples define the relationships between modules, thus resources' usages depend on the transferred tuples' size at a certain time. Total network usage is based on Eq. (24).

$$TNU = \frac{\sum_{i=1}^{N'} (L_i * S_i)}{T_{max}} \quad (24)$$

In Eq. (24), L_i and S_i are the latency and size of ith tuple overall, N' is the total number of tuples, and T_{max} is the maximum simulation time.

5.2.4 Application Delay

The delay of application execution is calculated by the system clock and the end time of a tuple.

$$T_{TN} = \begin{cases} SC - T_{st} & \text{if } T_a = 0, \\ \frac{T_{st} * N_{ET} + (SC - T_{st})}{N_{ET} + 1} & \text{if } T_a \geq 0 \end{cases} \quad (25)$$

The end time of the tuple is calculated by Eq. (25). Where T_{st} is the tuple start time, SC is the system clock, $(SC - T_{st})$ is the execution time, and N_{ET} is the number of executed tuple types. T_a is the average CPU time based on the tuple type. CC is the system clock, and ET is the emitting time of a tuple. T_s is transfer time between two modules.

$$T_{TR} = \frac{T_{st} * N_{RT} + SC - T_s}{N_{ET} + 1} \quad (26)$$

The tuple receipt time is based on Eq. (26). N_{RT} is the number of receipt tuple types. Application delay is calculated by the difference time between tuple end time in a module and tuple receipt time in another module.

5.3. Fairness

We evaluate the fairness of the offloading method based on the Jain index [32], which is computed as:

$$JainIndex = \frac{(\sum_{j=1}^Q x_j)^2}{Q * \sum_{j=1}^Q x_j^2} \quad (27)$$

Q is the number of devices that contributed to the offloading, and x_j is the jth device's energy consumption. The Jain index is between $\frac{1}{Q}$ and 1; a better offloading method has a more Jain index.

5.4 Comparison scenarios

This section categorizes our scenarios to analyze the proposed approach and other algorithms. Table 7 shows four different scenarios. Where Scenario 1 is based on a different number of users. Scenario 2 is considered for four different module sizes. Scenario 3 is for comparing the methods based on four mobile types. Also,

we compare our proposed approach with others based on different intervals. The reason for using the values introduced in the diagram is according to the type of application. Since this application is introduced in the iFogsim emulator, so it comes with values by default. We tried to consider less and more of these parameters to get a good estimate in terms of scalability, number and type of mobile devices, and module size.

Table 7: Comparison scenarios

| No. | Description | Values |
|------------|------------------|-------------------------|
| Scenario 1 | Number of users | 1, 3, 7, 10 |
| Scenario 2 | Module size (MB) | 1000, 2000, 5000, 10000 |
| Scenario 3 | Mobile types | AB, AC, BC, ABC |
| Scenario 4 | Interval (ms) | 100, 200, 500, 1000 |

5.5 Scenario 1: Comparison of offloading performance based on the number of users

One of the parameters to show the performance of the offloading methods is the number of users. Here, we compare the energy consumption, total execution cost, network usage, and delay of MUCAO in MEC by the number of users. As can be seen in Figs. 3, 4, 5, 6, and 7, there are values of the number of users by 1, 3, 7, and 10 in the horizontal axis.

Fog devices serve multiple users simultaneously. On the other hand, given the number of resources these devices have, when the number of users reaches a certain size, they reach a degree of optimization. This means that devices can perform resource management operations more efficiently. Due to the hierarchical structure of the network and users' distribution, it will be possible to improve the results even with the increase of users, which can be seen in the results.

5.5.1 Energy consumption based on the number of users

The analysis of Fig. 3 shows that the energy consumption of MUCAO is less than [the original and offload methods](#). MUCAO can decrease energy consumption in a higher number of users. As this figure, the maximum energy consumption is on the number of users by 7 for the original method by $1.635 * 10^7$ MJ. The minimum value is on the number of users by 10 for the FLUCO method by $1.58 * 10^7$ MJ. This result shows that the FL-based method with distributed structure causes less energy consumption than others. Also, adding context-awareness information to FLO and using more than one ED cause to create a better method as FLUCO for energy efficiency. [The reason for the improvement is a distributed algorithm of FLUCO that executes in multi EDs. The FLUCO causes MDs with lower resources to transfer more of their modules to EDs. As a result, the workload in the network has been distributed in a balanced way. There is not much difference between the energy consumption of methods with increasing the number of users. Since FLUCO distribute modules in the network, and also the capacity of devices is restricted. Thus the number of modules in devices cannot increase. Finally, we have not got more computations to calculate energy consumption for them.](#)

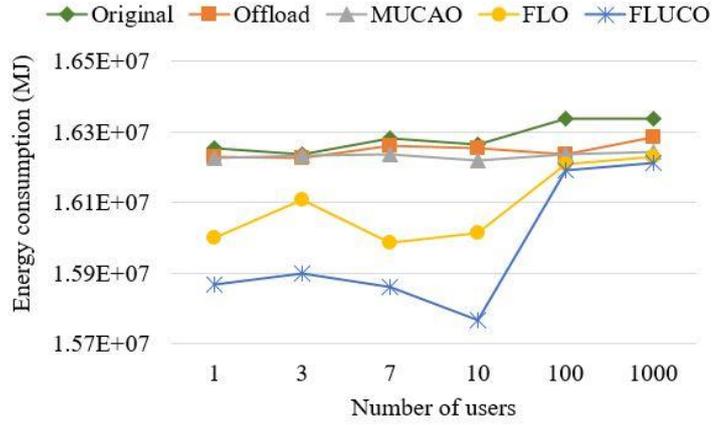


Fig. 3: Energy consumption based on the number of users

5.5.2 Total execution cost based on the number of users

The cost is one of the important metrics in this work. We can see in Fig. 4 with increasing the number of users from 3 to 7, and 10 causes increasing cost in original and offload methods. MUCAO has a balanced cost than these methods in many users with fluctuating between 4.12×10^6 \$ and 4.15×10^6 \$. Additionally, FLUCO with minimum energy consumption less than FLO and MUCAO is placed in the first rank of total execution cost. This result shows that with the increase in the number of users and distribution in the environment, the FLUCO has managed the cost well and brings economic savings. [The main reason for this improvement is our distributed algorithm in some EDs that cause choosing the best device with high performance and minimum delay.](#)

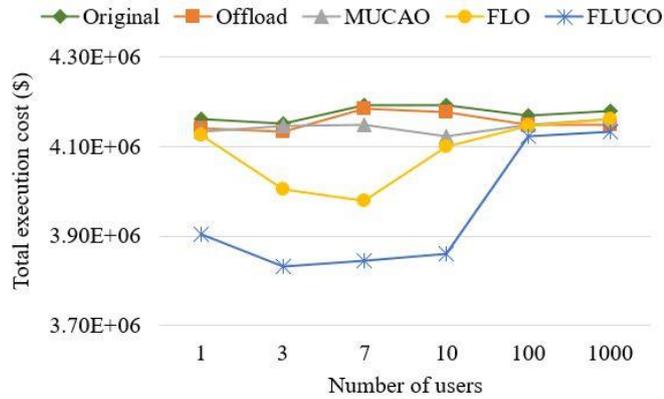


Fig. 4: Total execution cost based on the number of users

5.5.3 Network usage based on the number of users

Fig. 5 shows that MUCA than original and offload methods can increase network usage by considering context-awareness. As this figure, the higher number of users could not increase this metric so much. The reason for this result is the best matching of offload destination instead of first matching. MUCAO has used network resources better with higher performance than these two methods. On the other hand, by being superior to others, FLUCO and FLO managed resources better. [Since the main idea of these methods is distribution; as a result, the modules can offload to a wide range of devices. That is why all devices in the](#)

network are almost busy with minimum free time. FLUCO and FLO can be an obstacle for wasting time in devices.

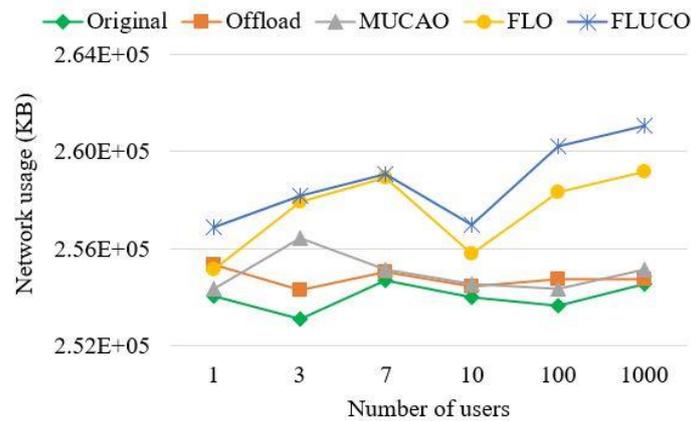


Fig. 5: Network usage based on the number of users

5.5.4 Delay based on the number of users

The delay of the application loop by MUCAO has a slight decrease except in users 3. These results are better than original and offload methods. As Fig. 6, the maximum delay is related to the original method in the number of users by 7. The minimum value is related to FLUCO in the number of users by 10. This showed that using context information and distributed algorithms cause to fast executing of requests and offloading process. This means MDs using FLUCO can quickly find the best destination to offload their modules and save more time.

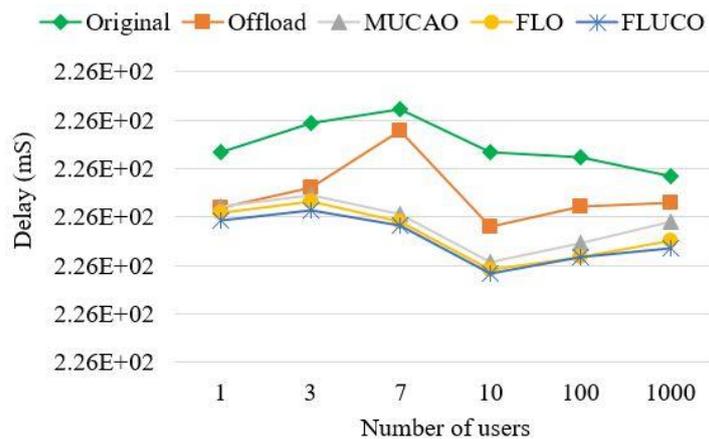


Fig. 6: Delay based on the number of users

5.5.5 Jain index based on the number of users

We provide the fairness of offloading algorithms by Jain index value in Fig. 7. Since the original method has not any offloading thus, we compare others. As we mentioned, this metric uses the energy consumption and the number of edge servers that contribute to offloading. The maximum Jain index is related to FLUCO in the number of users by 10. This shows in FLUCO; more edge servers are used to offload modules. Also, this result proves better load balancing in the FLUCO than others. As we mentioned in Eq. (27), energy

consumption is an important parameter here. The results also show that devices' proposed approaches cause to devices have a good cooperative for offloading modules from MDs to the best devices.

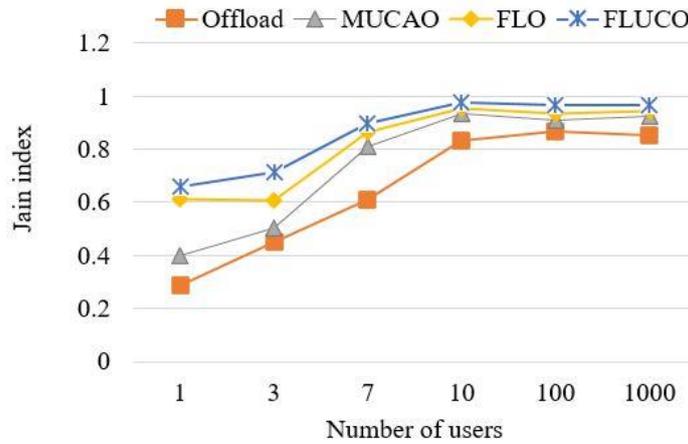


Fig. 7: Jain index based on the number of users

5.6 Scenario 2: Comparison of offloading performance based on module size

The module size is another metric for evaluating the offloading methods in this work. Based on Figs. 8, 9, 10, 11, and 12, there are module size's values by 1000, 2000, 5000, and 10000 in the horizontal axis.

5.6.1 Energy consumption based on module size

According to Fig. 8, the module size hasn't more effect on energy consumption. The original, offload, and MUCAO show almost equal energy consumption. On the other hand, FLUCO and FLO have got better results. This means using distributed structure and context awareness can improve the energy consumption of the system. Since EDs and cloud resources have more capacities than MDs, different modules can be offloaded and executed to the upper layer in the network. Also, regarding more devices and widely distributed modules, the modules size has not got a considerable change in energy consumption.

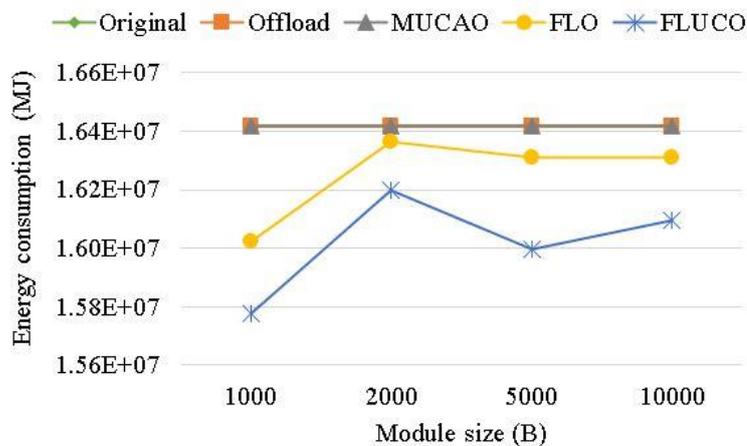


Fig. 8: Energy consumption based on module size

5.6.2 Total execution cost based on module size

Since different module sizes have no changes in energy consumption for original, offload, and MUCAO, as Fig. 9, the execution cost is without changes in different module sizes. The context-awareness of application, devices, and network environment allows MDs to have more selections for offloading their modules. As a result, the execution cost that is based on MIPS of devices calculates in a wide range. As a result, this has not got a lot of effects totally. However, FLUCO with minimum execution cost is the best method in comparison with others.

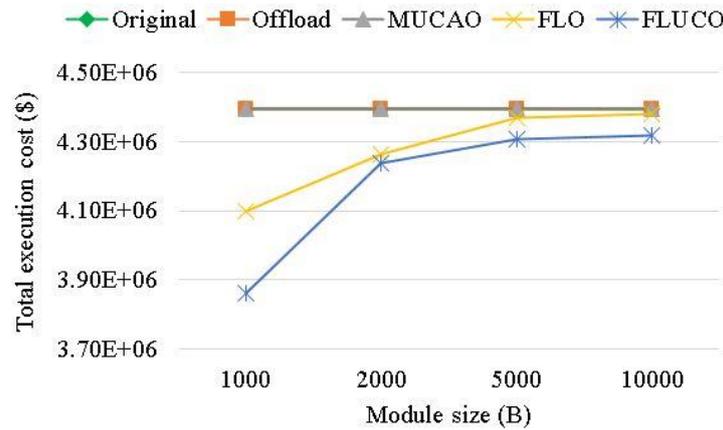


Fig. 9: Total execution cost based on module size

5.6.3 Network usage based on module size

Fig. 10 shows that the minimum network usage is related to the original method in 2000 B by 1.125×10^5 . FLUCO places 105 MB and its maximum value in 1000 B and network usage by 1.15×10^5 . Increasing module size causes increasing network usage. Also, in comparison with FLO, our proposed method as FLUCO, can improve network size. We need to explain increasing the module size causes raising total network usage, but this happens until a specific module size because EDs and cloud capacity is more than MDs. Thus, total network usage will not have many changes.

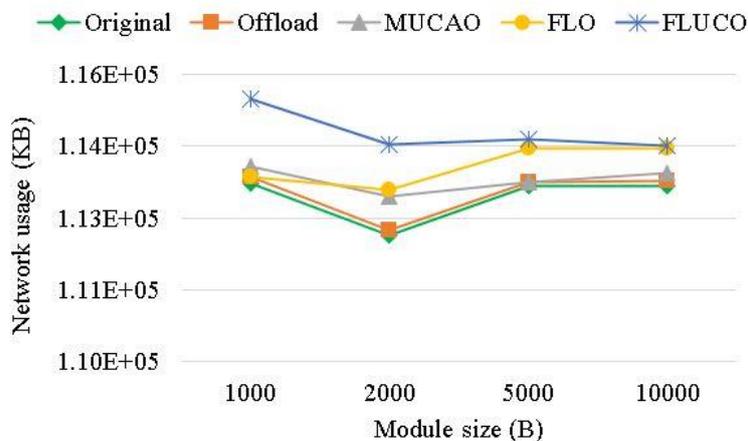


Fig. 10: Network usage based on module size

5.6.4 Delay based on module size

As shown in Fig. 11, increasing the module size causes a decrease in the delay of methods. The reason for decreasing delay in the original method is the local execution of modules. Also, we should consider some of the modules might not execute locally for not being enough resources. Of course, by this way, the energy consumption of MDs will be increase. Since edge servers have more capacity than module sizes, they can execute offloaded modules in less time. Also, FLUCO has better results than others in module sizes by 2000 and 5000 B. In module size 1000 and 10000 B, FLUCO has got a little improvement than FLO. This shows that a distributed algorithm can manage and offload them to the best devices when module size increases.

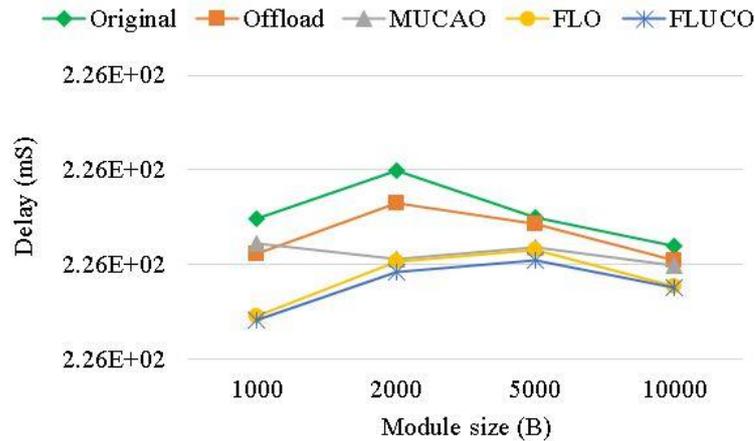


Fig. 11: Delay based on module size

5.6.5 Jain index based on module size

In Fig. 12, the fairness of the offloading method is between 0.6 and 0.8. However, the range of this metric in MUCAO is between 0.8 and 1.0. This proves MUCAO with considering context information is fairer than the offloading method. On the other hand, using distributes algorithms, and more EDs convert FLUCO to the best algorithm, among others. This means the energy consumption of all devices in the network was in a distributed manner. Also, decreasing the delay of modules cause all devices to consume less energy, so that the proposed approaches are better in this case.

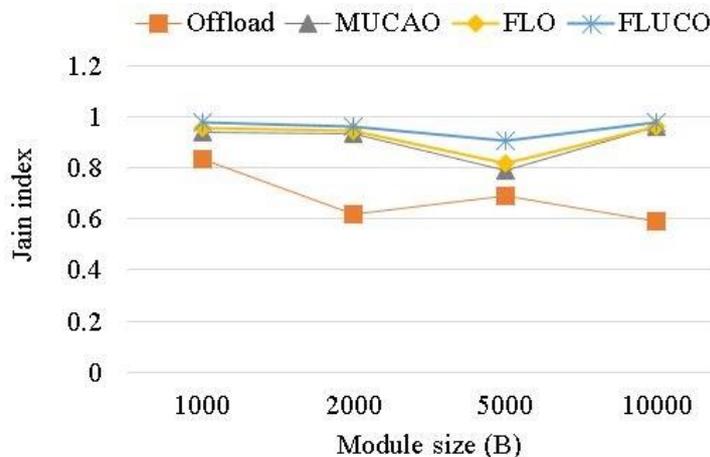


Fig. 12: Jain index based on module size

5.7 Scenario 3: Comparison of offloading performance based on mobile types

Figs. 13, 14, 15, 16, 17 are based on different mobile types, as shown in Table 6.

5.7.1 Energy consumption based on mobile types

Analysis of energy consumption based on mobile types shows better results of FLUCO than others in all states AB, AC, BC, ABC. The results in Fig. 13 proves the diversity of mobiles can cause less energy consumption by the proposed method. FLUCO, with minimum energy consumption of about 1.58×10^7 MJ, is the best method than others. This method shows that a distributed algorithm in different devices can offload modules with less energy consumption. Also, heterogeneous devices with different configurations have got required resources for local computation.

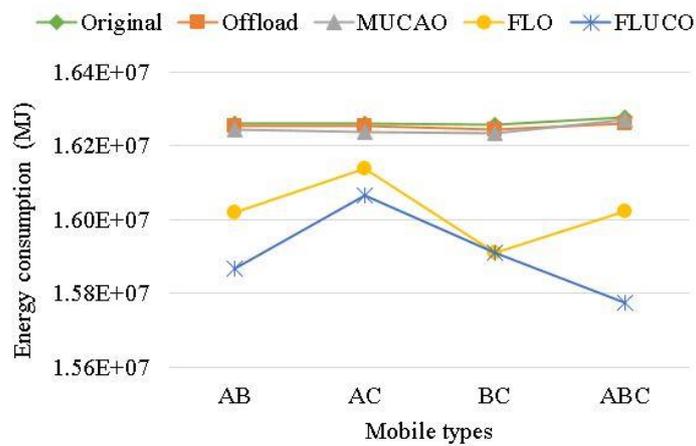


Fig. 13: Energy consumption based on mobile types

5.7.2 Total execution cost based on mobile types

According to Fig. 14, using different mobile types decreases total cost in all methods. FLUCO has better results than others. More capacity of CPU, memory, and battery causes mobile devices to execute more modules locally. This process decreases the offloading cost. Another reason for improving the FLUCO can be fair offloading in a wide range of devices. Thus, decreasing the cost in a device and distributing fair offloading to other devices can present lower cost.

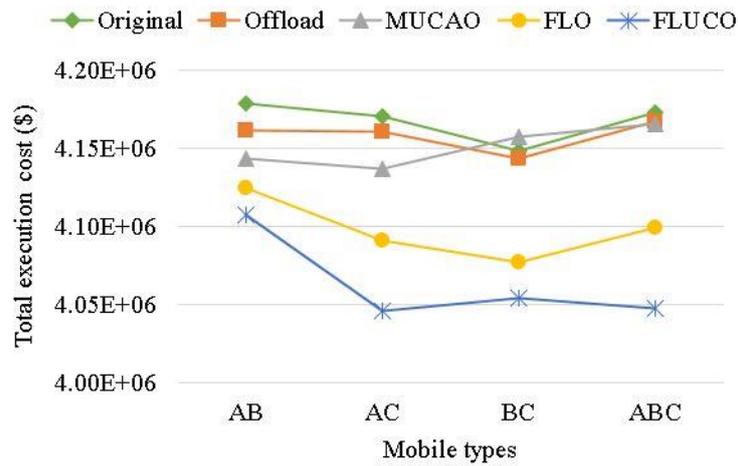


Fig. 14: Total execution cost based on module size

5.7.3 Network usage based on module types

Based on Fig. 15, network usage has a gradual increase by different mobile types. We can see in this figure that the FLUCO has maximum network usage in mobile type BC by 2.6×10^5 KB. Thus, diversity in mobile types has a direct effect on network usage. Using a distributed structure of the network causes more network usage in the MEC. This means the FLUCO with a distributed method can use many devices in the network and the number of jobless devices will decrease.

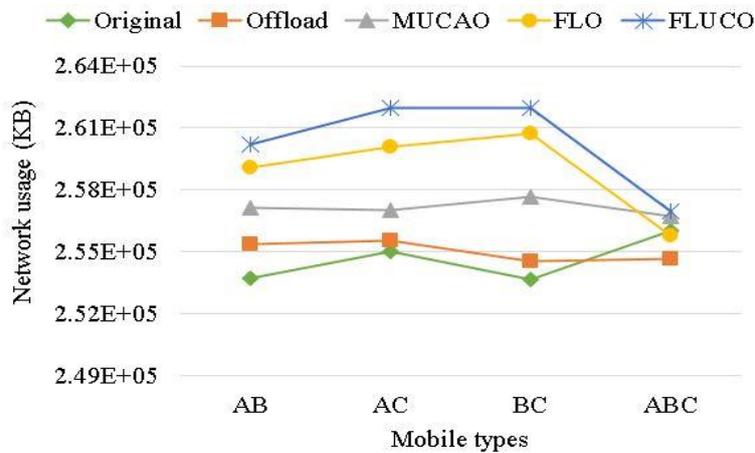


Fig. 15: Network usage based on module types

5.7.4 Delay based on module types

Fig. 16 shows that delay in all mobile types AB, AC, and BC has sensitive changes. FLUCO has a minimum delay equal to 226.2 mS on mobile type ABC. The results prove that the increase in the diversity of MDs causes more challenges in delay. The distributed algorithms as FLO and FLUCO can do better than others. FLUCO has got less delay when the mobile type is AC. The reason for this improvement can be the context awareness in FLUCO than FLO.

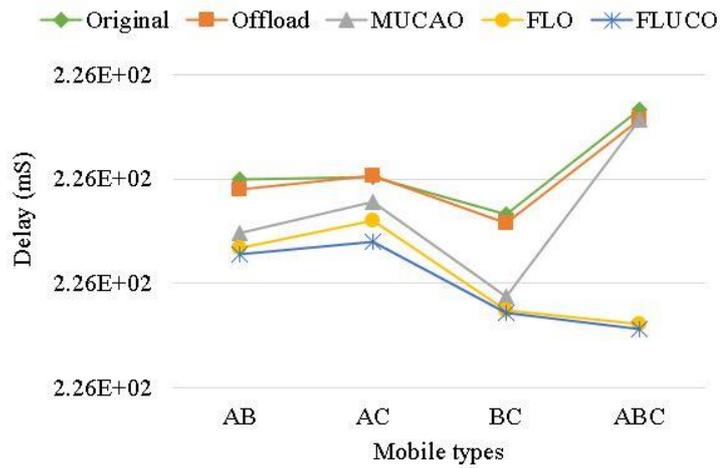


Fig. 16: Delay based on module types

5.7.5 Jain index based on module types

The fairness metric shows that FLUCO has the best results in all mobile types. Also, Fig. 17 proves the minimum fairness is related to the offloading method in mobile type AB by 0.6. Thus, distributed multi-user context-aware is a suitable offloading method in MEC. In fact, if the offloading method can place the modules on a wide range of devices, we will have a fairway.

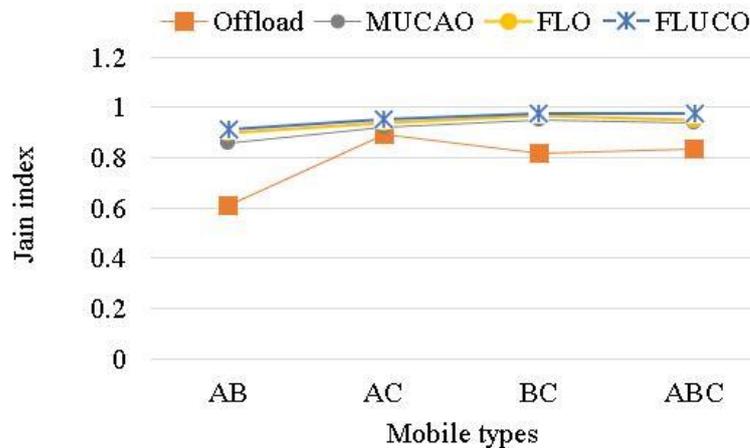


Fig. 17: Jain index based on module types

5.8 Scenario 4: Comparison of offloading performance based on interval

Offloading's interval shows the time distance between the resource management process. We control the workflow to the system by interval value. The energy consumption, total execution cost, network usage, delay, and Jain index are evaluated by offloading's intervals equal to 100, 200, 500, and 1000 mS. We set these values for the spacing between the input data goes back to the type of application. Since the application is intended to process input data in an average of 200 mS. Therefore, we consider numbers in the same range.

5.8.1 Energy consumption based on interval

Fig. 18 proves that our proposed MUCAO and FLUCO methods can decrease energy consumption than the original, offload, and FLO methods. FLUCO with 1.57×10^7 MJ is the best than others. Thus, FLUCO is very

suitable for offloading in the VRGAMEFOG application. Since the interval value means the distance time between the resource management process, increasing that causes the offloading method will have more time for process or offloading modules to best devices. As a result, we can see; generally, the FLUCO used this chance better than others.

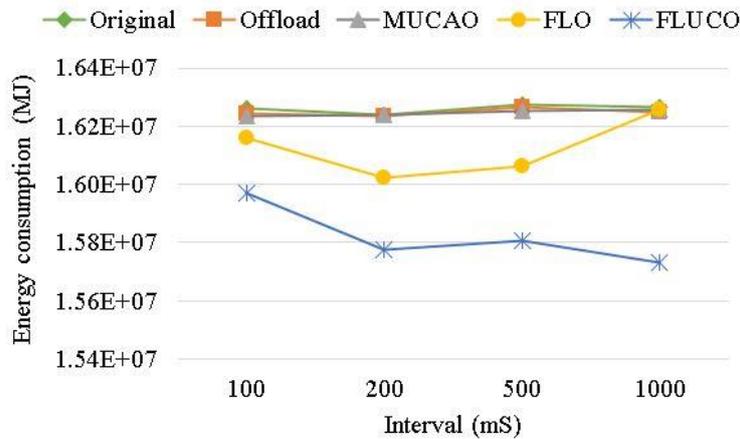


Fig. 18: Energy consumption vs interval

5.8.2 Total execution cost based on the interval

Analysis of execution cost by all offloading methods shows that interval equals 500 causes more energy consumption. According to Fig. 19, the minimum execution cost of 4.9×10^6 \$ is related to FLUCO by interval 200 mS. The worst execution is related to the original method by 4.18×10^6 \$ in the interval of 500 mS without any offloading. The results prove the superiority of the distributed algorithm over others in MEC. The fair using of resources in devices causes less cost, so FLUCO with the capability of distribution and context awareness has better results than others.

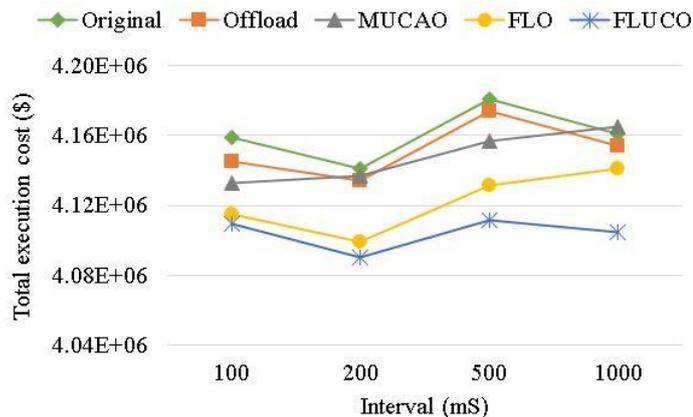


Fig. 19: Total execution cost vs interval

5.8.3 Network usage based on interval

The results of the simulation show the competition between all methods. They cause to near network usage values in the interval by 100, 500, and 1000 mS. Fig. 20 indicates, in the average stats, FLUCO is the best

offloading method than others. The main reason for this improvement is the distributed structure of FLUCO, also using the properties of devices, application modules, and environment are important in context-awareness.

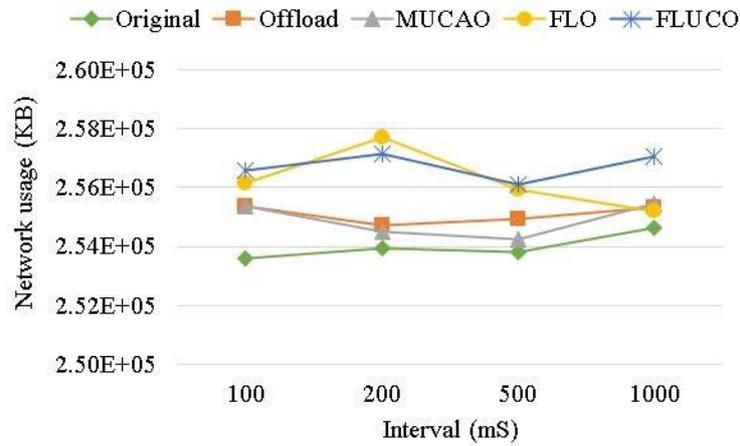


Fig. 20: Network usage vs interval

5.8.4 Delay based on the interval

Based on Fig. 21, the MUCAO method has a gradual increase in the delay parameter when the interval is grown. Of course, this MUCAO is an excellent method in intervals by 100 and 200 mS. According to this figure, the original method has the worst result in the interval by 1000 mS. The offload method has a fluctuated result with the lowest in the interval of 200 mS and higher in the interval by 500 mS. More analysis shows that FLO and FLUCO have got lower delays than others. The results show that these two methods can quickly offload the modules to the best devices with minimum delay. Of course, the lowest delay equal 226.3 mS, is related to FLUCO in an interval of 200 mS.

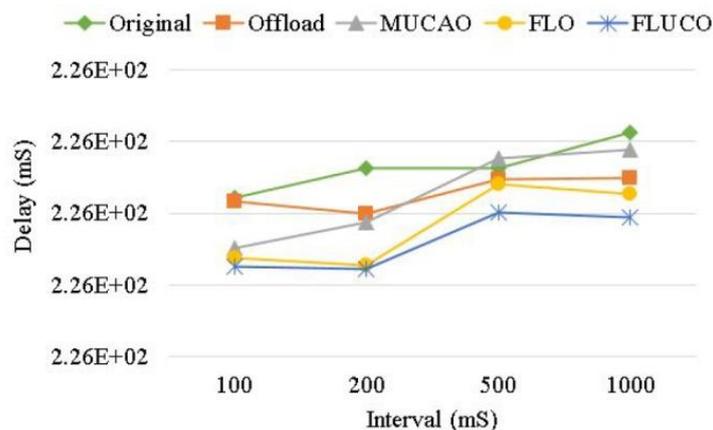


Fig. 21: Delay vs interval

5.8.5 Jain index based on interval

According to Fig. 22, MUCAO and FLUCO have a gradual increase in the Jain index with maximum fairness in interval of 1000 mS. However, the offloading method has fluctuated values, and it could

approximately close to MUCAO in 1000 mS. However, FLUCO with the highest Jain index is better than others. This proves that a distributed algorithm can be a fair offloading method. Thus, the dynamic context-awareness and distributed structure of the proposed algorithm can improve the performance of MEC.

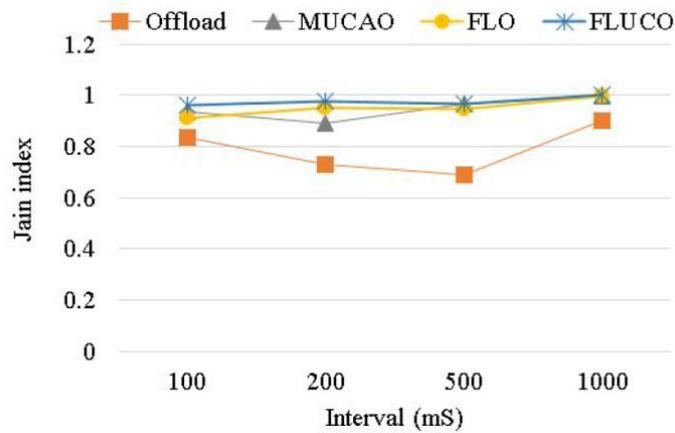


Fig. 22: Delay vs interval

6 Conclusion

Computation offloading in MEC is faced with many challenges. In this paper, we investigate context-aware offloading by considering multi-user. We also present a distributed algorithm as FLUCO to get close to the MEC structure. To solve this problem, a MAPE loop is used in all offloading processes. Our method helps MDs to offload their modules to edge servers or cloud if they can not execute those locally with less cost. The results show that FLUCO is superior to original, offload, MUCAO, and FLO methods in energy consumption by 2%, 2%, 2.1%, and 0.7% in total execution cost by 3%, 3%, 2.34%, and 1.08% network usage by 2%, 2%, 1.21%, and 0.001% delay by 0.01%, 0.01%, 0.005%, and 0.001% and 0.002%, respectively. Also, FLUCO is fairer than offload, MUCAO, and FLO methods in the Jain index by 18%, 4.01%, and 1.6%, respectively. These results prove that our proposed offloading algorithm with context-aware information and distributed structure could improve the network performance in the mentioned metrics. As future work, we work on MEC with FL-based methods on other case studies. Cooperative mobile crowding is another challenge of FL in MEC for more research. FL is vulnerable to communication security issues such as Distributed Denial-of-Service (DoS) and jamming attacks. Also, we will research on the protection of data privacy for MEC users.

References

1. Gourisaria, M. K., Samanta, A., Saha, A., Patra, S. S., & Khilar, P. M. (2020). An Extensive Review on Cloud Computing. In *Data Engineering and Communication Technology* (pp. 53-78). Springer, doi: 10.1007/978-981-15-1097-7-6.
2. Mutlag, A. A., Ghani, M. K. A., Arunkumar, N. A., Mohammed, M. A., & Mohd, O. (2019). Enabling technologies for Fog Computing in Healthcare IoT Systems. *Future Generation Computer System*, 90, 62-78. doi:10.1016/j.future.2018.07.049.
3. Ghobaei-Arani, M., Souiri, A., & Rahmanian, A. A. (2019). Resource management approaches in fog computing: a comprehensive review. *Journal of Grid Computing*, 1-42. doi:10.1007/s10723-019-09491-1

4. Wang, F., Xu, J., & Cui, S. (2020). Optimal Energy Allocation and Task Offloading Policy for Wireless Powered Mobile Edge Computing Systems. *IEEE Transactions on Wireless Communications*.
doi:10.1109/TWC.2020.2964765.
5. Hu, Y. Patel, M. Sabella, D. Sprecher, N. Young, V. (2015). Mobile edge computing a key technology towards 5G. *ETSI White Paper 11(11)(2015) 1-16*.
6. Mach, P., & Becvar, Z. (2017). Mobile edge computing: A survey on architecture and computation offloading. *IEEE Communications Surveys & Tutorials*, 19(3), 1628-1656.
doi:10.1109/COMST.2017.2682318.
7. Aral, A., Brandic, I., Uriarte, R. B., De Nicola, R., & Scoca, V. (2019). Addressing Application Latency Requirements through Edge Scheduling. *Journal of Grid Computing*, 17(4), 677-698.
doi:10.1007/s10723-019-09493-z
8. Rahman, A. U., Malik, A. W., Sati, V., Chopra, A., & Ravana, S. D. (2020). Context-aware opportunistic computing in vehicle-to-vehicle networks. *Vehicular Communications*, 24, 100236.doi:10.1016/j.vehcom.2020.100236.
9. Liang, Z., Liu, Y., Lok, T. M., & Huang, K. (2019). Multi-user computation offloading and downloading for edge computing with virtualization.*IEEE Transactions on Wireless Communications*,18(9),4298-4311.
doi:10.1109/TWC.2019.2922613.
10. Orsini, G., Bade, D., & Lamersdorf, W. (2018). CloudAware: empowering contextaware self-adaptation for mobile applications. *Transactions on Emerging Telecommunications Technologies*, 29(4), e3210.
doi:10.1002/ett.3210.
11. Luo, C., Goncalves, J., Velloso, E.,&Kostakos, V. (2020). A Survey of Context Simulation for Testing Mobile Context-Aware Applications. *ACM Computing Surveys (CSUR)*, 53(1), 1-39.
doi:10.1145/3372788.
12. Lim,W.Y.B., Luong, N.C., Hoang, D.T., Jiao, Y., Liang, Y.C., Yang, Q., Niyato, D. and Miao, C. (2020). Federated learning in mobile edge networks:A comprehensive survey. *IEEE Communications Surveys & tutorials*.
doi:10.1109/COMST.20020.2986024
13. Ren, J., Zhang, D., He, S., Zhang, Y., & Li, T. (2019). A Survey on End-Edge- Cloud Orchestrated Network Computing Paradigms: Transparent Computing, Mobile Edge Computing, Fog Computing, and Cloudlet. *ACM Computing Surveys (CSUR)*, 52(6), 1-36.doi:10.1145/3362031.
14. Boukerche, A., Guan, S., Grande, R. E. D. (2019). Sustainable Offloading in Mobile Cloud Computing: Algorithmic Design and Implementation. *ACM Computing Surveys (CSUR)*, 52(1), 1-37.doi:10.1145/3286688.
15. Tang, L. and He, S. (2018). Multi-User Computation Offloading in Mobile Edge Computing: A Behavioral Perspective, *IEEE Network*, 32(1), 48-53.
doi:10.1109/MNET.2018.1700119.
- 16 . Tran, D.H., Tran, N.H., Pham, C., Kazmi, S.M.A., Huh, E.-N., Hong, C.S. (2017). OaaS: offload as a service in fog networks. *IEEE Computing* 99(11), 1081–1104.doi:10.1007/s00607-017-0551-z.
17. Zhan,W., Luo, C., Min, G.,Wang, C., Zhu, Q.,&Duan, H. (2020). Mobility-Aware Multi-User Offloading

- Optimization for Mobile Edge Computing. *IEEE Transactions on Vehicular Technology* 69(3)3341-3356.
doi:10.1109/TVT.2020.2966500.
18. Kuang, L., Gong, T., OuYang, S., Gao, H., & Deng, S. (2020). Offloading decision methods for multiple users with structured tasks in edge computing for smart cities. *Future Generation Computer Systems*.
doi:10.1016/j.future.2019.12.039.
19. Kuang, Z., Shi, Y., Guo, S., Dan, J., & Xiao, B. (2019). Multi-User Offloading Game Strategy in OFDMA Mobile Cloud Computing System. *IEEE Transactions on Vehicular Technology*, 68(12), 12190-12201.
doi:10.1109/TVT.2019.2944742.
20. Yang, X., Fei, Z., Zheng, J., Zhang, N., & Anpalagan, A. (2019). Joint Multi- User Computation Offloading and Data Caching for Hybrid Mobile Cloud/Edge Computing. *IEEE Transactions on Vehicular Technology*, 68(11), 11018-11030.doi:10.1109/TVT.2019.2942334.
21. Liu, Z. Z., Sheng, Q. Z., Xu, X., Chu, D., & Zhang, W. E. (2019). Context-aware and Adaptive QoS Prediction for Mobile Edge Computing Services. *IEEE Transactions on Services Computing*.doi:10.1109/TSC.2019.2944596.
22. Lin, W., Peng, G., Bian, X., Xu, S., Chang, V., & Li, Y. (2019). Scheduling Algorithms for Heterogeneous Cloud Environment: Main Resource Load Balancing Algorithm and Time Balancing Algorithm. *Journal of Grid Computing*, 17(4), 699-726.
doi:10.1007/s10723-019-09499-7
23. Chang, Z., Zhou, Z., Ristaniemi, T., Niu, Z. (2017). Energy efficient optimization for computation offloading in fog computing system. In: *GLOBECOM 2017–2017 IEEE Global Communications Conference*, pp.1–6.
doi:10.1109/GLOCOM.2017.8254207.
24. Liu, L., Chang, Z., Guo, X., Mao, S., Ristaniemi, T. (2018). Multi-objective optimization for computation offloading in fog computing. *IEEE Internet Things J.* 5(1), 283–294.
doi:10.1109/JIOT.2017.2780236.
25. Al-Ayyoub, M. Al-Quraan, M. Lo'ai A. (2017). Delay-aware power optimization model for mobile edge computing systems. Springer, Volume 21, pp 1067–1077.doi:10.1007/S00779-017-1032-2.
26. Huang, L., Feng, X., Zhang, L., Qian, L., & Wu, Y. (2019). Multi-server multi-user multi-task computation offloading for mobile edge computing networks. *Sensors*, 19(6), 1446.
doi:10.3390/s19061446.
27. Salehan, A., Deldari, H. & Abrishami, S. (2019) An online context-aware mechanism for computation offloading in ubiquitous and mobile cloud environments. *J Supercomput* 75, 3769–3809 .
doi:10.1007/s11227-019-02743-7.
28. Cho, J., Sundaresan, K., Mahindra, R., Van der Merwe, J., & Rangarajan, S. (2016, December). ACACIA: context-aware edge computing for continuous interactive applications over mobile networks. In *Proceedings of the 12th International on Conference on emerging Networking EXperiments and Technologies* (pp. 375-389).
doi:10.1145/2999572.2999604.

29. Nawrocki, P., & Sniezynski, B. (2017). Autonomous context-based service optimization in mobile cloud computing. *Journal of Grid computing*, 15(3), 343-356.doi:10.1007/s10723-017-9406-2
30. Baraki, H., Jahl, A., Jakob, S., Schwarzbach, C., Fax, M., & Geihs, K. (2019). Optimizing Applications for Mobile Cloud Computing Through MOCCAA. *Journal of Grid Computing*, 17(4), 651-676.doi:10.1007/s10723-019-09492-0
31. Chen, X., Chen, S., Zeng, X., Zhang, Y., Zheng, X., & Rong, C. (2017). Framework for context-aware computation offloading in mobile cloud computing. *Journal of Cloud Computing: Advances, Systems and Applications*, 6(1), 1–17.doi:10.1186/s13677-016-0071-y.
32. Ghasemi-Falavarjani, S. Nematbakhsh, M. Ghahfarokhi, B. S. (2015). Contextaware multi-objective resource allocation in mobile cloud, *Computers & Electrical Engineering*, vol. 44, pp. 218-240.doi:10.1016/j.comeleceng.2015.02.006.
33. Roostaei, R. Movahedi, Z. (2016). Mobility-aware and fault-tolerant computation offloading for Mobile Cloud Computing. *Intl IEEE Conferences on Ubiquitous Intelligence & Computing, Advanced and Trusted Computing*.
34. Kosta, S. Aucinas, A. Hui, P. Mortier, R. Zhang, X. (2012). Thinkair: dynamic resource allocation and parallel execution in the cloud for mobile code offloading. In *IEEE Proceedings of INFOCOM*, pp. 945–953.
doi:10.1109/INFOCOM.2012.6195845.
35. Lin, T., Lin, T. A. Hsu, CH. King, CH. (2013). Context-aware decision engine for mobile cloud offloading. *IEEE Wireless Communications and Networking Conference Workshops (WCNCW)*.111-116.
doi:10.1109/WWCNCW.2013.6533324.
36. Kotb, Y., Al Ridhawi, I., Aloqaily, M., Baker, T., Jararweh, Y., & Tawfik, H. (2019). Cloud-based multi-agent cooperation for IoT devices using workflow-nets. *Journal of Grid Computing*, 17(4), 625-650.
doi:10.1007/s10723-019-09485-z
37. Ren, J., Wang, H., Hou, T., Zheng, S., & Tang, C. (2019). Federated learning-based computation offloading optimization in edge computing-supported internet of things. *IEEE Access*, 7, 69194-69201.
DOI:10.1109/ACCESS.2019.2919736.
38. Yang, K., Jiang, T., Shi, Y., & Ding, Z. (2020). Federated learning via over-the-air computation. *IEEE Transactions on Wireless Communications*, 19(3), 2022-2035.doi:10.1109/TWC.2019.2961673.
39. Wang, X., Han, Y., Wang, C., Zhao, Q., Chen, X., & Chen, M. (2019). In-edge ai: Intelligentizing mobile edge computing, caching and communication by federated learning. *IEEE Network*, 33(5), 156-165.
doi:10.1109/MNET.2019.1800286.
40. Shen, S., Han, Y., Wang, X., & Wang, Y. (2019). Computation Offloading with Multiple Agents in Edge-Computing-Supported IoT. *ACM Transactions on Sensor Networks (TOSN)*, 16(1), 1-27.
doi:10.1145/3372025.
41. A. Computing, et al., An architectural blueprint for autonomic computing ,IBM White Paper 31(2006) 1-6.
42. Gupta, H., Vahid Dastjerdi, A., Ghosh, S. K., Buyya, R. (2017). iFogSim: A toolkit for modeling and simulation of resource management techniques in the Internet of Things, Edge and Fog computing environments. *Software: Practice and Experience*, 47(9), 1275-1296.doi:10.1002/spe.2509.

43. Sutton, R. S and Barto, A. G. Reinforcement learning: An introduction. MIT press, 2018.
44. Hua, P. Wu-Shao, W. Ming-Lang, T. Ling-Ling, L. (2020). Joint optimization method for task scheduling time and energy consumption in mobile cloud computing environment. *Applied Soft Computing Journal* 80:534–545.
doi: 10.1016/j.asoc.2019.04.027
45. Peng, K. Zhu, M. Zhang, Y. Liu, L. Zhang, J. Leun, V. C. M. and Zheng, L. (2019). An energy- and cost-aware computation offloading method for workflow applications in mobile edge computing. *EURASIP Journal on Wireless Communications and Networking* 207:1-15.
doi: 10.1186/s13638-019-1526-x
46. Nawrocki, P. Sniezynski, B. (2020). Adaptive Context-Aware Energy Optimization for Services on Mobile Devices with Use of Machine Learning. *Wireless Personal Communications* 115:1839–1867.
doi: 10.1007/s11277-020-07657-9
47. Burd, T. D. Brodersen, R. W. (1996). Processor Design for Portable Systems. *Journal of VLSI Signal Processing Systems* 13, 203-221.
doi: 10.1007/BF01130406