

# Time-optimized sequential decision making for service management in smart city environments

Saleh ALFahad<sup>a,\*</sup>, Christos Anagnostopoulos<sup>b</sup> and Kostas Kolomvatsos<sup>c</sup>

<sup>a</sup> School of Computing Science, University of Glasgow, Glasgow, UK

ORCID: <https://orcid.org/0000-0003-0101-4458>

<sup>b</sup> School of Computing Science, University of Glasgow, Glasgow, UK

ORCID: <https://orcid.org/0000-0003-1517-6757>

<sup>c</sup> Department of Informatics & Telecommunications, University of Thessaly, Volos, Greece

ORCID: <https://orcid.org/0000-0002-9442-3340>

Received 9 December 2022

Accepted 10 February 2023

**Abstract.** Edge Computing is a new computing paradigm that aims to enhance the Quality of Service (QoS) of applications running close to end users. However, edge nodes can only host a subset of all the available services and collected data due to their limited storage and processing capacity. As a result, the management of edge nodes faces multiple challenges. One significant challenge is the management of the services present at the edge nodes especially when the demand for them may change over time. The execution of services is requested by incoming tasks, however, services may be absent on an edge node, which is not so rare in real edge environments, e.g., in a smart cities setting. Therefore, edge nodes should deal with the timely and wisely decision on whether to perform a service replication (*pull-action*) or tasks offloading (*push-action*) to peer nodes when the requested services are not locally present. In this paper, we address this decision-making challenge by introducing an intelligent mechanism formulated upon the principles of optimal stopping theory and applying our time-optimized scheme in different scenarios of services management. A performance evaluation that includes two different models and a comparative assessment that includes one model are provided found in the respective literature to expose the behavior and the advantages of our approach which is the OST. Our methodology (OST) showcases the achieved optimized decisions given specific objective functions over services demand as demonstrated by our experimental results.

**Keywords:** Service management, sequential decision making, task offloading, edge computing, smart cities, optimal stopping theory

## 1. Introduction

As a result of the fast growth of the Internet of Things (IoT), smart cities are replete with linked sensors and computing devices used for information collection and processing [42]. Various collaborative services (e.g., smart

---

\*Corresponding author: Saleh ALFahad, Knowledge and Data Engineering Systems, School of Computing Science, University of Glasgow, Glasgow, UK. E-mail: [2426473A@student.gla.ac.uk](mailto:2426473A@student.gla.ac.uk).

medical assessment, smart university, elderly support) emerge as a result of the right processing and utilization of the huge data created by the Internet of Things [45]. On the basis of integrated services, it is likely that the smart city will be able to increase the Quality of Experience (QoE) of people while reducing the use of resources. Furthermore, smart cities are endowed with a robust capacity to observe, analyze, and adapt to their residents [34]. Nevertheless, because of the rapid growth of data and the limited processing capacity of endpoints, it is challenging to develop joint services that are typically time and latency-sensitive when depending purely on IoT devices. To increase the QoE in smart cities, it is essential to use additional computational resources [13]. As a result, Edge Computing (EC), which completely leverages computational resources at the network edge and can meet real-time demands, has been proposed as an effective paradigm to resolving the issue of inadequate computing resources in smart cities [57].

In order to offer ubiquitous and pervasive services, the present expansion of IoT and EC makes them possible to have a large number of devices and processing nodes located close to end-users [30]. In recent years, we have also seen the introduction of several mobile devices and applications, such as drones and Vehicular Networks (VN), each with its own unique set of capabilities. The development of this technology has made it possible for computing and sensing devices to launch applications such as Augmented Reality (AR), predictive analytics activities at the edge, intelligent vehicle control, traffic management, and interactive applications [4]. These types of services are given in response to the need for processing data being gathered by IoT devices and, then, being sent to EC nodes [63]. A range of IoT applications may offer to end users with more precise and detailed network services [18]. In this situation, an increasing number of equipment and sensors are being networked through IoT technology, and this equipment and sensors will create vast amounts of data that need additional processing, therefore delivering service providers and end users with intelligence. In the traditional Cloud computing setup, all information must be transferred to centralized servers, and after processing, the outcomes must be sent directly to the requestors, e.g., sensors or other devices. This procedure places significant stress on the network, particularly in terms of bandwidth and resource requirements for data transmission. Furthermore, when data size increases, bandwidth utilization will decline.

EC has been suggested to address the aforementioned difficulties [24,59]. The fundamental goal is to increase the processing capabilities at the network edge; specifically, computing [51], bandwidth [62], and storage resources [43] are relocated nearer to IoT systems in order to minimize core data traffic and response delay and to support resource-intensive IoT applications. In comparison to Cloud [11,25], EC nodes have limited storage and computing capabilities. Nonetheless, they are able to host a number of services that make it easier to carry out a variety of processing tasks [40]. These kinds of activities are often presented in the form of assignments that ‘request’ the processing of data (e.g., predictive analytics, explanatory-driven models). The fact that IoT and EC may work together to facilitate the implementation of collaborative activities in which nodes can collaboratively complete the tasks that have been requested is an intriguing development [22]. For instance, nodes share services or data and they even offload processing duties onto peers in certain circumstances. Because of their reduced capability for storage and computing, EC nodes are only able to host a (sub-)set of the total accessible services and the data that has been obtained. The collected data can be reported by stream monitoring e.g., the evolution of a certain phenomenon [6,27,28] and streams generated by autonomous nodes (e.g., unmanned vehicles) [29]. Upon these streams, we can support the extraction of knowledge, the detection of events or any other processing [30].

Services are essential in order to perform the processing responsibilities that have been delegated by applications or users in the form of tasks. As a result, the demand for services is always adjusting to accommodate the dynamic requests of various activities/tasks. This indicates that a task can be asking for a service that may or may not be available locally. Then, nodes have two options:

- Option 1: They can perform a service replication (*pull-action*) from their edge neighbors or the Cloud, but only if the service is suitable for their computing capabilities;
- Option 2: They can delegate the responsibility (*push-action*) to the peers/Cloud that presently host the service(s).

The aforementioned options deal with deciding where to transfer service in order to keep the Quality of Service (QoS) at a high level. Because of the nature of the EC environment and the non-static behavior of end users, it is obviously difficult to provide an ideal migration/replication plan that can be implemented. Moreover, due to the situations that make the local authority of the user complex, offloading tasks involves sending them to peers or the

Cloud for processing. For instance, the node that is assigned the duty of processing as the receptor of a task could have a high load, and as a result, it might not be able to guarantee an effective completion within the permitted time interval that may be specified by the requester. A strategy for the migration of services was presented in [60]; its purpose is to satisfy the service latency requirements of EC while keeping migration costs and trip times to a minimum. Furthermore, in [20], the authors recommended using a reinforcement learning-based model to solve the service migration issue.

In this paper, we focus on the pull-action leaving the initiative to EC nodes to enhance their autonomous nature and suggest that services' replication may be optimized and controlled in the EC environment by adopting the principles of optimal stopping theory (OST) [21]. In this case, EC nodes must decide locally when it is the *best* time to replicate the service in order to maximize its ability to carry out the requested tasks. We have to mention that the decision for services replication is made by every EC node *independently* being fully aligned with the needs of the incoming tasks and the status of the node. The remaining paper is organized as follows. In Section 2, we provide a summary of the prior work as well as the rationale and our contribution, while Section 3 delves into the specifics of the proposed OST-based decision-making model. The outcomes of our performance and comparative assessment are presented in Section 4, and Section 5 draws the conclusion of this work and indicates potential avenues for further research.

## 2. Related work & contribution

### 2.1. Related work

The majority of the approaches for task offloading at the edge depends on the selection of whether or not the task should be processed locally or should be offloaded to peers or Cloud. The time in execution latency [26] and the amount of energy used should both be reduced as much as possible. These are the two primary goals. The work presented in [35] provides a framework for offloading computation workloads from a user device to a server hosted in EC. Our approach is different from that due to the fact that our time-optimized decision primarily considers the actual benefit of payoff which indicates the current users' requests for specific service. In [35], the User Equipment (UE) choice on whether to offload computing tasks with the highest CPU availability for a certain application is driven by the radio network information service (RNIS), according to the current value of round trip time (RTT) between the UE and the edge node. In our method, the offloading choice ideally relies on the number of users' requests for specific service as well as the cost of delay. Therefore, the decision of offloading is made when the dedicated time of the required service by users is finished which indicates the required service is not worth being hosted. As a result, this will optimize the edge users' capacity.

The work in [9] presents a method for optimizing the decision-making process that an Unmanned Aerial Vehicle (UAV) uses to choose whether or not to do the compute task locally or to offload it to an edge server. The decision is made based on a series of interactions that take place between the UAV and IoT systems. During these interactions, the UAV receives feedback on the state of the network and EC server, which enables an estimation of the remaining amount of time needed to complete the task. As a result, the UAV uses this information to solve an optimization problem with the purpose of decreasing a weighted sum of delay and energy costs as much as possible. In our work, EC chooses the ideal moment that has a high payoff to perform a service replication (pull-action) from their edge neighbors or the cloud. The approach described in [1] analyses the present state of the node's resources, which are modeled as a Markov Decision Process utilizing Q-learning for training, in order to identify which portion of the application ought to be offloaded and then moves on with an offloading decision. The primary objective is to reduce the amount of time that the offloaded applications are delayed while taking into account the mobile fog that is located nearby, the mobile fog that is located adjacent, or the Cloud as suitable offloading locations. In our work, we have examined the scenario in which the edge node only performs a service replication (pull-action) from its edge neighbors or the cloud since the edge node is static and makes this decision in a short delay as close to the real's delay scenario.

The work in [17] focuses on reducing the processing delay of all activities and the energy usage of the edge device by optimizing the task allocation decision and the central processing unit (CPU) frequency of the edge

device. However, this study focuses on making decisions based on the popularity of the required service to complete the task. The authors of [16,31] studied the offloading problem decision for a multi-edge user & multi-edge node. However, in our approach, we focus on making the decision of (pull or push-action) for a single user to edge node or peer since our approach OST has strong deals with low capacity better than the approach in [16,31]. Using a time-division multiple access protocol in [50], edge users can conduct their individual activities locally or offload all or a portion of them to the Access Point (AP). Their objective is to decrease the AP's overall energy consumption relative to the user at the network's edge. Nevertheless, our recommended solution is subject to a choice depending on the service's popularity and our stopping time with the maximum payoff. The study presented in [12] addresses the issue of task offloading in ultra-dense networks with the goal of reducing latency as much as possible while preserving the battery life of user equipment. On the other hand, our method provides the optimum option in terms of (pull or push-action) the service (from or to) the edge's neighbors or peers.

The work in [23,41] primarily formulate the decision-making of offloading based on a resource scheduling factor. However, in our approach, we develop the decision of (push or pull-action) w.r.t. accumulation of the newly received tasks and cost incurring. A state-action-payoff-state-action (RL-SARSA) method is introduced in [3] to tackle the resource management issue at the edge server and make the best offloading option for reducing system cost. Moreover, the work in [49] aims to make the offload decision based on minimizing the cost. However, our work emphasizes making the decision with the greatest payoff, considering the expense cost each time that we do not make the decision of (push-action). In the study presented in [10,48], the issue of offloading computation for many users in a wireless environment with uncertainty is investigated. Nevertheless, the goal of our work is to provide the single-edge user with the best possible time in which to make a decision (push or pull-action). In [32,55], the authors examine the challenge of partial offloading scheduling and resource allocation for edge computing systems with various independent activities. The objective is to reduce the weighted total of processing latency and energy usage while satisfying the tasks' transmission power limitations. Nevertheless, our approach enables the edge user to choose the optimal (pull-action) solution. Since the service has previously been hosted, the wait for future tasks will be reduced. The work in [15,38] investigate energy-efficient task offloading in edge computing. The goal is to reduce the amount of energy required for task offloading. However, our method focuses on making the choice at the optimal moment to host the most popular tasks, hence optimizing the edge's capacity resource utilization.

Migration of data and services may be used to fill 'gaps' in the local knowledge base of a processing node. In [54], the interested reader may obtain a survey of related initiatives. When migrating services, it will be efficient to meet processing demands locally; nonetheless, a list of factors must be considered. To minimize overburdening the network, service migration may be accomplished in phases; service components should be "hooked up" to the hosting infrastructure swiftly, easily, and automatically. One of the technologies used for these reasons is Machine Learning (ML). ML may serve as the foundation for producing models that predict the movement of users and provide a proactive procedure for services migration [8]. Reinforcement learning may improve agent behavior while determining the optimal action to take [2,52]. The difficulty in depending on supervised machine learning technology is identifying the appropriate representative training data set. As a result of the raised degree of uncertainty in the corresponding decision-making, it is challenging to gather data that includes all alternative node statuses. The service migration model may alternatively be expressed as an online queue stability issue [37]. The issue may be addressed using either a Lyapunov optimization or a multi-objective optimization scheme [47,56], and the Pareto optimum solution can be determined. In the past, efforts have mostly focused on minimizing delay within the restrictions of energy usage in EC settings. However, the bulk of research efforts that aim at a user-centric service migration model is 'affected' by the extra time required to address the optimization issue, except if a list of assumptions is made or approximate solutions can be accepted. Using Markov chains is also a possible solution to settle the issue. Using a modified policy-iteration algorithm, the authors of [53] provide a solution to a finite-state Markov decision process. A technique based on Thompson sampling is suggested in [36] to facilitate dynamic service migration choices. The authors of [46] describe a service allocation module for networked automobiles powered by Cloud-based services. The objective is to increase the amount of autonomy of automobiles capable of picking nearby automobiles to share adopted services. The authors of [26] suggest a proactive system for determining the effective administration of services and tasks that are present/reported at EC nodes. The suggested model monitors the demand for existing services and rationalizes their management, i.e., their local presence/invocation when the demand is modified by the

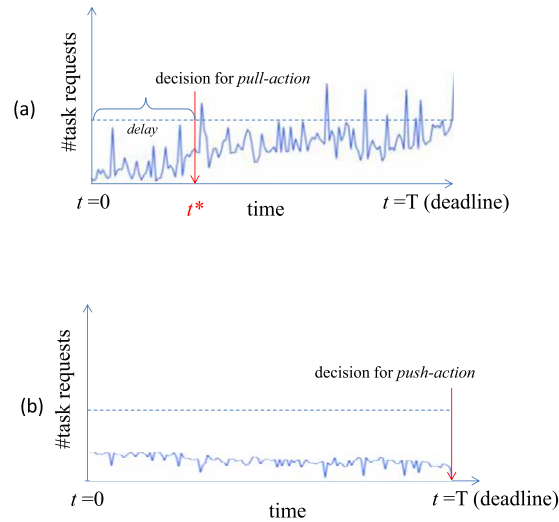


Fig. 1. Observing the task requests rate at each time instance  $t$  to decide on (a) a high popularity or (b) low popularity of the service and act intelligently based on the history of the task requests.

desired processing operations. To proactively achieve the strategic objectives of the envisioned model, a statistical inference process is initiated for each incoming task.

## 2.2. Rationale & contribution

Figure 1 illustrates in a simplified way the problem of deciding whether to perform a service replication (pull-action) from their edge neighbors or the cloud or offload (push-action) it to the cloud service to make the necessary computation to complete the required task. Let us assume that the starting time of recording the number of requests for a task is  $t = 0$  and the end time period, hereinafter referred to as deadline is  $T > 0$ , where a decision should be made. Figure 1 shows two decisions against time. Specifically, in Fig. 1 (a), we observe the scenario of making the pull-action decision at an early time  $t^*$  (and before the deadline  $T$ ) since we are expecting receiving more tasks requesting a specific service. This denotes that there is a relatively high demand for this service, thus, it is advisable to decide on pulling the service locally to the node. In Fig. 1 (b), we observe the scenario in which we delayed the decision making in light of increasing our confidence that there is no high demand of the service. Once the deadline elapses and we have not decided on a pull-action, which indicates that the service request rate was relatively low, then, the node decides on offloading the service requests (push-action) either to its peers or the Cloud. In other words, this implies that there was not a large volume of tasks requesting the specific service, thus, offloading these tasks is necessary. Assume the case where we made the decision to download/pull the service earlier enough, and then after a while, we realized that the service was requested only by a small number of users. We would have consumed additional cost of resources without benefiting from the download/pull service. As a result, it would have been more efficient to delay our decision until the end of the dedicated time and then perform the action of offloading as it is shown in Fig. 1 (b). On the contrary, if we had decided to send/offload the request to the Cloud at an early time, then, we would not have high confidence whether this request was highly requested on the node, thus, this decision would not be efficient. The reason is that this service would have been highly requested in the future, thus, the push-action would not be appropriate. However, if after a while we realize that the service demand has significantly increased, then an early push-action results in overloading the network with offloading tasks to other peers or the Cloud unnecessarily. Thus, a proper process is to sophisticatedly delay any decision and perform service replication (pull-action) from node's neighbors or the cloud to the node at a specific time  $t^*$  as shown in Fig. 1 (a). Hence, the problem has been identified as finding that time  $t^*$  to ensure that our resources are utilized in the best way trading off between task offloading and service replication. A fundamental solution has been proposed, which ensures the appropriate decision is taken at the time  $t^*$  taking into account the delay costs.

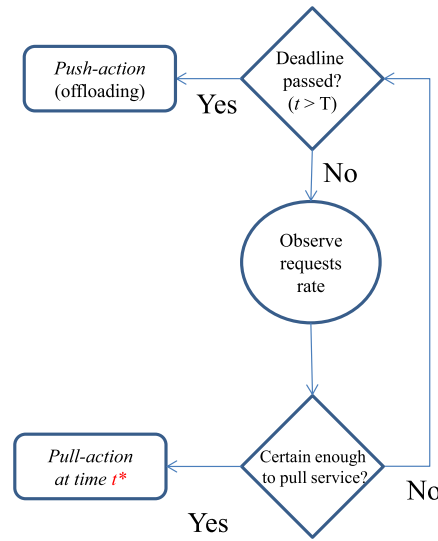


Fig. 2. Sequential decision-making diagram on either task offloading (push-action) or service replication (pull-action) at the best time  $t^* < T$  by the edge node.

**Contribution.** Our approach is adaptable to applications that use offloading decision-making algorithms in EC environments. Our major technical contribution is:

- We provide a time-optimized intelligent method that enhances the likelihood of making the decision of service replication (pull-action) of the service with the greatest payoff. This is reflected by the cumulative sum of the number of task requests to the node under consideration taking into account the incurred cost.
- We provide theoretical analysis of the optimality of our model based on the principles of the OST.
- We report on a detailed comparative evaluation of our OST-based method against the Ideal Decision ruLe (IDL), which produces the actual maximum payoff, and a heuristic Deterministic Rule (DR).
- A comparative assessment is provided comparing our results with the related work in [5], which is based on the Secretary Optimal Stopping Time problem.

In addition, Fig. 2 highlights the state-space of the decision-making. Specifically, Fig. 2 describes the mechanism of the seeking the best time for making the decision: either service replication (pull-action) or task offloading (push-action). Mainly, the decision is affected by the number of the received tasks at each time instance  $t$  for a specific service.

The rationale has as follows. The node receives a number of incoming task requests  $Y_t$  at every time instance  $t < T$ . Then, the node has two options:

- If the number of the received requests is high enough, therefore the decision will be to pull the service at a time  $t \leq t^* < T$ .
- Otherwise, the node reconsiders its decision at the next time  $t + 1$ , where new service requests are coming.

If the deadline  $T$  elapses, then the decision will be to push the task to a peer or the Cloud for further processing. This is the first study that we are aware of that sequential decision making, which considers the decision of pull/push-action as an OST problem in EC environments. The overall goal is to increase the possibility of making the decision of service replication (pull-action) with the greatest payoff. Moreover, since the nodes are located among the servers that have been distributed at the edge of the network, our approach will be applied to EC applications such as Virtual Networks (VN) [61], UAV [4], and data mining applications [30].



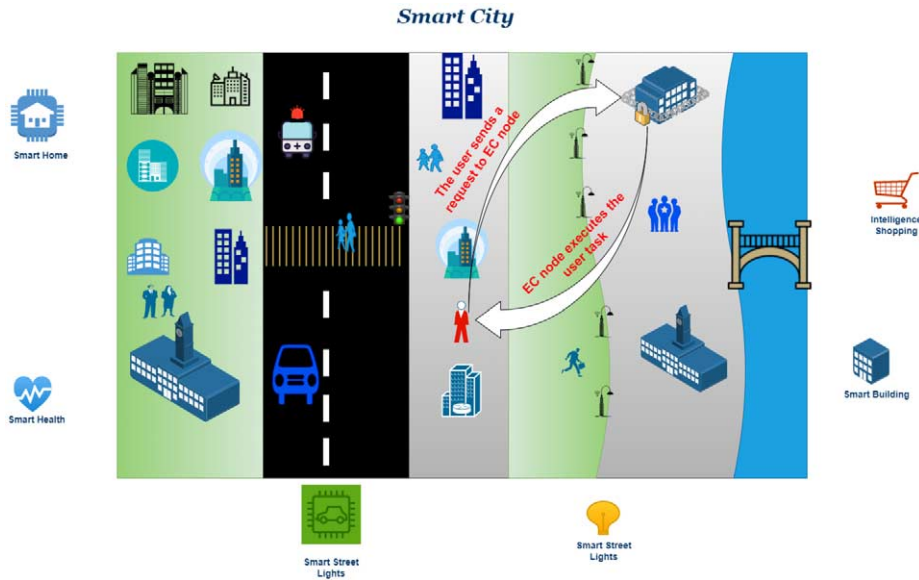


Fig. 3. Example of the two actions: pull and push in a smart city environment.

### 3. System model & problem formulation

In Section 3.1, we begin by providing an overview of the system model, and then in Section 3.2 we discuss the problem definition. In Section 3.3, the optimal stopping rule problem is introduced. Finally, in Section 3.4, we will give full details about cost-based sequential decision-making.

#### 3.1. System model

As seen in Fig. 3, a smart city is increasingly establishing vast IoT networks with widely distributed IoT devices that provide vast amounts of services. In light of the vast size and extensively spread nature of IoT networks, EC has become a strong and efficient paradigm for providing processing capabilities to IoT devices at the network's edge [58]. In EC, IoT services may be executed on EC, which offers low delay and reduces the bandwidth load. Nevertheless, it remains difficult to enhance the entire EC execution performance, for instance, resource utilization, EC energy usage, and load balance. However, we believe that the environment of the network is made up of a global cloud data center and an EC system [14,33]. The EC system is constructed of  $M$  base stations as defined in Table 1, and since each one is supplied with some EC servers, they collectively constitute an EC node. In accordance with it as well, we will refer to the set of EC nodes as  $\mathcal{M} = [1, 2, \dots, M]$ . Our general assumption is that multiple service providers would deploy their own applications on each EC node. Each edge device has the ability to delegate tasks to an EC node that is equipped with appropriate computing services. In addition, edge devices have the capability of off-loading tasks to a global cloud, presuming that the global cloud has all of the required services and a decent amount of computing resources [19]. Edge devices make the decision of task service replication (pull-action) in accordance with the maximum payoff value. The time is segmented into frames denoted by the notation  $t = 1, 2, \dots, T$ . A decision of service replication (pull-action) is determined by the edge device at some moment in time  $t^*$ , depending on local information (e.g., task information like the number of services that have been received). We rely on the assumption that the location of the edge device and the network environment remain the same throughout each time instance. This allows us to guarantee that the task of replication (pull-action) decision will be effective at the same time instance  $t$ .

Table 1  
Notations of parameters

NOTATIONS	Definition
$\mathcal{M}$	Set of EC nodes.
$T$	Deadline.
$t$	Time instance.
$t^*$	Optimal stopping time.
$Y_t$	Number of tasks received at time $t$ .
$\lambda$	Arrival rate of tasks $Y$ .
$\theta$	Tolerance threshold.
$C$	Cost of the decision delay.
$R_t$	Accumulation of number of tasks $Y_t$ .
$\alpha$	Heuristic factor in DR model.
$H_t$	Payoff at time instance $t$ .

### 3.2. Problem definition

Within the scope of the task model, we rely on the assumption that the formation of the user's tasks follows a Poisson distribution with  $Y_t$ : the number of the task's requests at time instance  $t$  with the rate  $\lambda$  [44]. Due to the fact that the Edge device is constant, the number of EC nodes that the edge device is capable of sensing and receiving tasks during each time instance  $t$ . As a result, the cumulation of the number of requests up to  $t$  is applied to every task that are received by the edge device. The decision of service replication (pull-action) is subjected to  $\theta$  which is a number of tasks that our OST model is based on it to take the action of making the decision. The cumulative distribution function of the poison distribution can be expressed as in (1).

$$P(Y \leq y; \lambda) = \sum_{k=0}^y \frac{e^{-\lambda} \cdot \lambda^k}{k!} \quad (1)$$

Let us define the cumulation of the number of the requests  $R_t, 0 \leq t \leq T$  as:

$$R_t = \sum_{k=0}^t Y_k. \quad (2)$$

The objective is to get as near as possible to  $\theta$ . In the case that  $R_t > \theta$ , the observation is immediately terminated and the service offloading scenario is activated, since we have acquired much more usable information than requested. Clearly, the second scenario is not unwanted, but it results in a 'penalty' that we ought to have prevented before  $R_t > \theta$ . This penalty is measured by the formula  $R_t > \theta$ . However, we may include a penalty for the increased latency caused by receiving one more task before stopping. Each time we do not stop, hence not triggering the service replication (pull-action) scenario, we incur a cost greater than zero. This cumulative cost up to  $t$  requires the procedure to decide whether to stop or continue with another observation by evaluating the trade-off between extending the observation for potentially more usable data and the extra latency required to activate the service offloading scenario. This price may reflect the urgency of a monitoring application that demands a service offloading option in near real-time. Alternatively, if the application requires very precise service offloading option outcomes, a (limited) delay must be permitted.

Based on the above interpretation, we provide the following payoff function  $H_t(R_t)$ , which involves (i) the task threshold  $\theta$ , (ii) the delay cost per monitoring  $C \geq 0$ .

$$H_t(R_t) = \begin{cases} R_t - C \cdot t & \text{if } R_t \leq \theta \\ 0 & \text{if } R_t > \theta \end{cases} \quad (3)$$



The rationale behind this payoff is that we desire the rate of service demands per time  $t$  taking out the delay cost to be significantly increasing. This expresses the likelihood that the pull-action is stochastically a better decision than the push-action. This is because we increase our confidence that in the (near) future we are expecting more requests for a service in a node. Therefore, we aim at maximizing this payoff and, hence, the corresponding likelihood. Our goal is to attain  $\theta$  by maximizing the number of the received tasks. If  $R_t$  is greater than  $\theta$ , then the return is zero since the mechanism should have triggered the service offloading scenario. Therefore, any additional delay was of no benefit.

It is mandatory to have a base result to be compared with our fundamental approach (OST). Therefore, we define the IDL, which measures the accumulative sum of the received tasks from users to the node. However, each time instance  $t$  the expected cost  $C \cdot t$  must be in the account. As a result, the instance time  $t$  that has the maximum payoff that results from IDL, is the best time  $t^*$  instance for making the decision (pull or continue). Moreover, we can express the IDL best time instance  $t^*$  as (4):

$$t^* = \left\{ 0 \leq t < T : \max_t \sum_{k=0}^t Y_k - C \cdot k, \text{ s.t. } \sum_{k=0}^t Y_k \leq \theta \right\} \quad (4)$$

We suggest the deterministic DR. DR is mainly based on calculating the number of received tasks  $Y$  after passing the amount of time  $t$ . In this rule, the decision of offloading is made if  $R_t$  of the received task at time instance  $t$ , which is determined by  $\alpha \cdot \frac{\theta}{C}$ , is less than theta  $\theta$ . Also, It is important to pay attention to the cost per each time instance ( $C \cdot t$ ). Then, the decision of service replication (pull-action) is expressed in DR as in (5). The pseudo-code is provided in Algorithm 1.

$$t^* = \min \left\{ 0 \leq t < T : \sum_{k=0}^t Y_k > \theta \text{ or } t > \alpha \frac{\theta}{C} \right\}. \quad (5)$$

Moreover, our theory has been compared to the work in [5]. Their approach is based on choosing the stopping time using SECPRO model. The stopping time decision in [5] was based on choosing the highest value of the received requests per time instance:

$$Y^* = \max_t (Y_t), \quad (6)$$

---

**Algorithm 1:** Heuristic decision-making model (DR).

---

**Input:** Tolerance threshold  $\theta$ ; observation cost  $C > 0$ ; heuristic factor  $\alpha \in [0, 1]$

**Output:** decided stopping time  $t^*$ .  $Stop \leftarrow False$ ;  $t \leftarrow 0$ ;  $R_0 \leftarrow 0$ ;

**repeat**

**observe** number of requests  $Y_t$  ;

$R_t \leftarrow R_{t-1} + Y_t$  /\* update the cumulative sum of requests \*/;

**if** criterion in (5) is satisfied **then**

$Stop \leftarrow True$ ;

$t^* \leftarrow t$  /\* activate the service replication (pull-action) decision and start-off a new task's service\*/;

**else**

$t \leftarrow t + 1$  /\* continue with the next received task\*/;

**end**

**until**  $t \leq T$

**if** ( $Stop = False$ )

        The (push-action)

**end**

---

**Algorithm 2:** Secretary-based OST model (SECPRO).

---

**Input:** Tolerance threshold  $\theta$ ; observation cost  $C > 0$ .  
**Output:** Secretary problem's optimal stopping time  $t^*$ .  
 /\* the algorithm observes task requests in the duration  $[0, \frac{T}{\exp}]$  \*/;  
 Stop  $\leftarrow$  False;  $t \leftarrow 0$ ;  $R_0 \leftarrow 0$ ;  $Y^* \leftarrow 0$ ;  
**for**  $t = (0 : \frac{T}{e})$   
   **observe** number of requests  $Y_t$  ;  
   **if**  $Y^* < y_t$  **then**  
      $Y^* = y_t$   
   **end**  
    $R_t \leftarrow R_{t-1} + Y_t$  /\* update the cumulative sum of requests \*/;  
**end**  
 /\* then the algorithm observes requests in the interval  $t \in \{\frac{T}{\exp}, T\}$  \*/;  
**for**  $t = (\frac{T}{e} + 1 : T)$   
   **observe** number of requests  $Y_t$  ;  
   **if**  $R^t > \Theta$  **then**  
     Push-action  
   **end**  
    $R_t \leftarrow R_{t-1} + Y_t$  /\* update the cumulative sum of requests \*/;  
   **if** criterion in (7) is satisfied **then**  
     Stop  $\leftarrow$  True;  
      $t^* \leftarrow t$  /\* activate the service replication (pull-action) decision and start off a new task's service\*/;  
   **else**  
      $t \leftarrow t + 1$  /\* continue with the next received task\*/;  
   **end**  
**until**  $t \leq T$   
**if** (Stop = False)  
   The (push-action)  
**end**

---

by edge node during the period from  $t \in \{0, \frac{T}{e}\}$ . Furthermore, taking into account calculating the cumulative sum of the number of the received tasks ( $R_t$ ) and the cost  $C$ . Then, the node will resume receiving tasks during the period of  $\{\frac{T}{e}, T\}$  where the first value  $Y_t > Y^*$  is the preferred stopping time of this theory. If the first value is detected, it stops immediately. This time is considered the stopping time for this theory. Also, we would like to emphasize that during the period  $\{\frac{T}{e}, T\}$  the cumulative sum of the received tasks ( $R_t$ ) is calculated beside the value of the cost as shown in (7). The pseudo-code is provided in Algorithm 2.

$$t^* = \min \left\{ \frac{T}{e} \leq t \leq T : \text{such that } \sum_{k=0}^t Y_k \leq \theta \text{ and } Y_t > Y^* \right\} \quad (7)$$

As a result, we found surprising and noteworthy results that demonstrate the superiority of our theory OST. This is what we will explain later in Section 3.3.

### 3.3. Sequential decision making based on OST

#### 3.3.1. Optimal stopping rule problem

The concept of optimal stopping [7] addresses the issue of selecting a time instance to conduct a specific action. The measure is taken to minimize an anticipated loss (or maximize an expected payoff). A stopping rule problem

is characterized by: (i) a set of stochastic variables, whose joint distribution is presumed to be known, and (ii) a set of loss functions  $(S_t(y_1, \dots, y_t))_{1 \leq t}$  or payoff functions  $(H_t(y_1, \dots, y_t))_{1 \leq t}$  that rely exclusively on the observed values  $y_1, \dots, y_t$  of related random variables. An optimal stopping rule problem is expressed as follows: We are monitoring the series of  $(Y_t)_{1 \leq t}$ , and at each time instance  $t$ , we decide whether to stop monitoring or proceed. If we cease monitoring at instance  $t$ , we incur loss  $S_t$  or gain payoff  $H_t$ . We would like to pick a stopping rule or time that reduces our estimated loss or, equivalently, optimizes our expected payoff.

**Definition 1.** The optimal stopping time  $t^*$  minimises the incurred loss  $E[S_{t^*}] = \inf_{0 \leq t \leq T} E[S_t]$ . In the case of payoff,  $t^*$  maximises the expected payoff, i.e.,

$$E[H_{t^*}] = \sup_{0 \leq t \leq T} E[H_t]. \quad (8)$$

Note that  $T$  may be  $\infty$ . The information accessible up to time  $t$  is a series  $F_t$  of values of the random variables  $Y_1, \dots, Y_t$ , which is also known as filtration.

**Definition 2.** The 1-stage look-ahead (1-sla) stopping rule (time) is defined as:

$$t^* = \inf\{t \geq 0 : H_t \geq E[H_{t+1}|F_t]\}. \quad (9)$$

This means that  $t^*$  requires stopping at the first time instance  $t$  where the payoff  $H_t$  for stopping at  $t$  is as large as the anticipated payoff of proceeding to the next time instance  $t + 1$  and then stopping.

**Definition 3.** Let  $B_t$  represent the occurrence  $H_t \geq E[H_{t+1}|F_t]$ . A stopping rule problem is monotonous if  $B_0 \subset B_1 \subset \dots$  nearly certainly. The set  $B_t$  is the set for which the 1-sla rule requires pausing at time instance  $t$ . The condition  $B_t \subset B_{t+1}$  indicates that if the 1-sla rule requires pausing at time  $t$ , it will likewise require stopping at time  $t + 1$  regardless of the value of  $y_{t+1}$ . In a same manner,  $B_t \subset B_{t+1} \subset B_{t+2} \subset \dots$  states that if the 1-sla rule requires pausing at time  $t$ , it will also require stopping at all future times, regardless of the future observations.

**Theorem 1.** The 1-sla rule provides the best solution for monotonous stopping rule issues.

*Proof.* See [39]. □

The edge device is required to complete making the decision of service replication (pull-action) of the service of the number of tasks  $y_t$  that it has received. Because of this, the primary goal is to increase the payoff attached to the decision of service replication (pull-action). The edge node is actively monitoring a continuous sequence of number of tasks that have been received. These tasks are being accumulated with those that have been received in the past w.r.t. performance criteria, which is denoted by the current delay cost  $C$  per time  $t$ . At each received new tasks, the node should decide whether to make the decision of service replication (pull-action) or not. The challenge arises from the fact that the node wants to establish a service replication (pull-action) policy that will increase the chances of making the decision of service replication (pull-action) with a high payoff. As a result, the node needs to go for the optimal decision time that maximizes the payoff, which is globally rated as being the same as or near to the actual (ideal) payoff.

### 3.4. Cost-based sequential decision making

**Problem 1.** Given the tolerance threshold  $\theta$  and delay cost  $C > 0$ , determine the optimal stopping time  $t^*$  such that  $\sup_{1 \leq t \leq \infty} E[H_t]$  is attained. The greatest expected return is then  $E[H_{t^*}]$ .

The Problem 1 is resolved by establishing a model based on the OST. Consider that the starting time  $t = 0$  and the deadline is  $T$ , with  $t < T$ . Our OST model aims to stop receiving more tasks when the cumulative sum of the received task  $R_t$  is near to  $\theta$ , taking into account the cumulative delay cost up to that point  $t$ . We provide a 1-sla

stopping rule based on the optimal stopping principle stated in Theorem 1, whereby we stop at the first time  $t$  such that:

$$H_t(R_t) \geq E[H_{t+1}(R_{t+1})|F_t], \text{ with the event } \{R_t \leq \theta\} \in F_t. \quad (10)$$

That is, any further receiving tasks at time  $t + 1$  would not contribute to maximizing the payoff. When the difference  $E[H_{t+1}(R_{t+1})|F_t] - H_t(R_t)$  is constant non-increasing with  $R_t$ , the 1-sla rule is optimal.

**Theorem 2.** *Given a series of payoff random variables  $R_1, \dots, R_t$ , the optimal stopping rule  $t^*$  for Problem 1 is provided in (11):*

$$t^* = \inf \left\{ t \geq 0 : \sum_{y=0}^{\theta-R_t} (R_t + y) - C(t+1) \cdot P(Y = y) \leq R_t - C \cdot t \right\} \quad (11)$$

*Proof.* Given that  $R_t \leq \theta$ , the conditional expectation  $E[H_{t+1}(R_{t+1})|R_t \leq \theta]$  is given by:

$$\begin{aligned} E[H_{t+1}(R_{t+1})|R_t \leq \theta] &= E_Y[H_{t+1}(R_t + Y)|R_t \leq \theta, R_t + Y \leq \theta]P(R_t + Y \leq \theta) \\ &\quad + E_Y[H_{t+1}(R_t + Y)|R_t \leq \theta, R_t + Y > \theta]P(R_t + Y > \theta) \\ &= E_Y[(R_t + Y) - c(t+1)|Y \leq -R_t]P(Y \leq \theta - R_t) \\ &= \sum_{y=0}^{\theta-R_t} (R_t + y) - C(t+1) \cdot P(Y = y). \end{aligned}$$

Therefore, the mechanism terminates at the first time  $t$ , where  $R_t$  is such that  $E[H_{t+1}(R_{t+1})|R_t \leq \theta] \leq R_t - C \cdot t$ .  $\square$

The 1-sla rule is optimal, since the corresponding difference is constantly non-increasing with  $R_t$  when  $R_t < \theta$ . A high-level description of the provided Theorem 2 is that, it is always better to stop at the first time  $t$  rather than at  $t + 1$ , because the current pay off at that  $t$  will be bigger than the pay off at the next time instance  $t + 1$ , and any other time  $t + \tau$  with  $\tau > 1$ . This is the principle of 1-sla rule, where our problem falls into this category. Moreover, the difference of the pay off at the optimal time  $t$  and  $t + 1$  is always positive. However, that difference tends to decrease with the time  $t$ . This indicates that at the very first time  $t$  where the pay off is bigger than the future expected pay off, then it is always beneficial to stop at that  $t$  and not at any other future time, since the difference of the consecutive pay off values is always decreasing. Therefore, the 1-sla rule in Theorem 1 is optimum, OST with fixed  $\theta$  and taking into account the delay cost  $C$  may ensure that the anticipated cumulative sum of the received tasks based on the stopping criteria is as near as possible to  $\theta$ , however, no other stopping rule can make such a claim. Based on the cost  $C$ , OST is adaptable for handling and controlling the latency and freshness of the service replication (pull-action) and offloading decision's variables. OST is concerned with the delay cost of monitoring and assumes that all received tasks are current for the application that uses the service replication (pull-action) and offloading decision. In this case, the OST is forced to terminate receiving new tasks to prevent waiting for an extended period of time, particularly when  $\lambda$  is very small. Therefore a large number of number of tasks are necessary to sum up to  $\theta$ . The pseudo-code is provided in Algorithm 3.

## 4. Experimental evaluation

### 4.1. Experimental setup

In our experiments, we experiment with four different models: (i) The proposed OST model. (ii) The heuristic DR model, which depends on the heuristic factor  $\alpha$ . In DR, the decision of pull-action is taken if the accumulation  $R_t$  at time instance  $t$ , determined by  $\alpha \cdot (\theta/C)$ , is greater than  $\theta$  provided in (5). Otherwise, if the deadline elapses, the decision of offloading is made. (iii) The SECPRO approach introduced in [5]. SECPRO measures the current

**Algorithm 3:** Time-optimized sequential decision making (OST).

---

**Input:** Tolerance threshold  $\theta$ ; observation cost  $C > 0$ .  
**Output:** Optimal stopping time  $t^*$ .  
 $Stop \leftarrow False$ ;  $t \leftarrow 0$ ;  $R_0 \leftarrow 0$ ;  
**repeat**  
  **observe** number of requests  $Y_t$ ;  
   $R_t \leftarrow R_{t-1} + Y_t$  /\* update the cumulative sum of requests \*/;  
  **if** criterion in (11) is satisfied **then**  
     $Stop \leftarrow True$ ;  
     $t^* \leftarrow t$  /\* activate service replication (pull-action) and start-off a new task's service\*/;  
  **else**  
     $t \leftarrow t + 1$  /\* continue with the next received task\*/;  
  **end**  
**until**  $t \leq T$   
  **if** ( $Stop = False$ )  
    The (push-action)  
  **end**

---

Table 2  
Experimental parameter setting

Notation	Parameter	Value/range
$T$	Deadline	100.
$\lambda$	Arrival rate of tasks	{3, 10, 30}.
$\theta$	Tolerance threshold	{50, 100, 150, 200, 600, 1200, 1800, 2400}.
$C$	Delay cost	{0.1, 0.3, 0.5, 0.6, 0.9}.
$\alpha$	Heuristic factor	{0.01, 0.05, 0.07, 0.09, 0.1, 0.3, 0.5, 0.6, 0.9}.

payoff when the decision has been made using the OST criterion in (7). (iv) The IDL is the real expected payoff ( $H$ ) gained provided in (4). This is the ground truth used to compare all the methods.

In the performance evaluation, we investigate the performance of the models: OST, DR, and IDL. Also, we set the values of for each of  $\lambda$ ,  $\alpha$  and  $\theta$ . We assigned  $\lambda \in \{3, 5, 10\}$  and  $\theta \in \{50, 100, 200\}$ . Moreover, we performed the experiments into three parts which are divided as the first experiment with  $\lambda$  and  $\theta$  set to 3 and 50, respectively. Then, the second experiment is set up with  $\lambda = 5$  and  $\theta = 100$ . Finally,  $\lambda = 10$  and  $\theta = 200$  for the third experiment. Also, the cost  $C$  was given a range of values and set these values for the three experiments. In fact,  $C = (0.01, 0.05, 0.1, 0.3, 0.5, 0.9) \cdot \lambda$  are the cost values. In addition, The  $\alpha$  was a range of numbers applied for all experiments in (0.01, 0.05, 0.07, 0.09, 0.1, 0.3, 0.45, 0.6, 0.9) as shown in Table 2. Moreover, in order to get a more precise result, we did each experiment 1000 times and then we calculate the mean values for all the results.

In our comparative assessment, we performed the experiment over all the models, i.e., OST, DR, IDL, and SECPRO. We set the values of  $\lambda$ ,  $\alpha$  and  $\theta$ . We assigned  $\lambda \in \{3, 30\}$ ,  $\alpha = 0.1$  and  $\theta \in \{50, 100, 150, 200\}$ ,  $\theta \in \{600, 1200, 1800, 2400\}$ . Moreover, we performed the experiment in two stages. The fourth experiment, which is the first stage, was set up with  $\lambda = 3$ ,  $\alpha = 0.1$ , and  $\theta \in \{50, 100, 150, 200\}$ . Furthermore, the second stage, which is the fifth experiment, was set up with  $\lambda = 30$ ,  $\alpha = 0.1$ , and  $\theta \in \{600, 1200, 1800, 2400\}$  as shown in Table 2. We run the experiments 1000 times and took the average values for all the obtained results.

## 4.2. Performance evaluation

### 4.2.1. Experiment 1

The settings for the Experiment 1 are:  $\lambda = 3$  and  $\theta = 50$  (low service demand rate). Here, we performed the experiment with the defined  $\lambda$  and  $\theta$ . Our experiment is mainly affected by the cost  $C$  and  $\alpha$ . The experiment

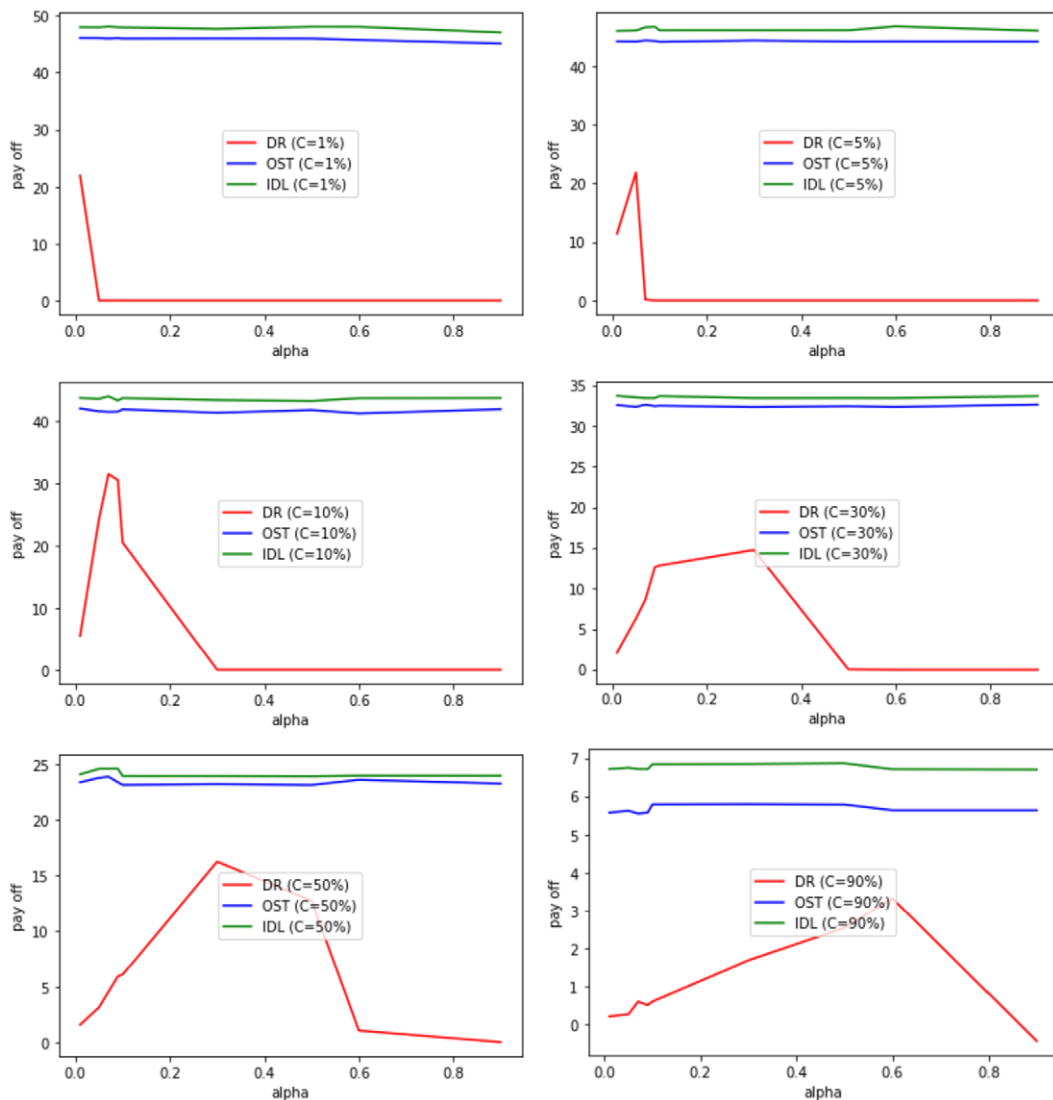


Fig. 4. (Experiment 1) The expected payoff vs delay cost with low service demand rate  $\lambda = 3$  and tolerance  $\theta = 50$  for OST, DR and IDL models.

measures the payoff and the time that decision has been made. Moreover, the experiment is performed over three models. The first model is IDL which is being chosen to be the baseline to be compared with the other two models. The second model is our fundamental approach which is OST. Finally, the third model is DR. The experiment test three different models which are the real Figure 4 shows the different effects of the cost all over the three models and obviously shows the result of payoff and how our approach is as close as IDL compared with DR.

#### 4.2.2. Experiment 2

The settings for the Experiment 2 are:  $\lambda = 10$  and  $\theta = 200$  (medium service demand rate). Here, we performed the experiment with the defined  $\lambda$  and  $\theta$ . Our experiment is mainly affected by the cost  $C$  and  $\alpha$ . The experiment measures the payoff and the time that decision has been made. Moreover, the experiment is performed over three models. The first model is the IDL model which is being chosen to be the baseline to be compared with the other two models. The second model is our fundamental approach which is OST. Finally, the third model is the DR model. Figure 5 shows the different effects of the cost all over the three models and obviously shows the result of payoff and how our approach is as close as IDL compared with DR.



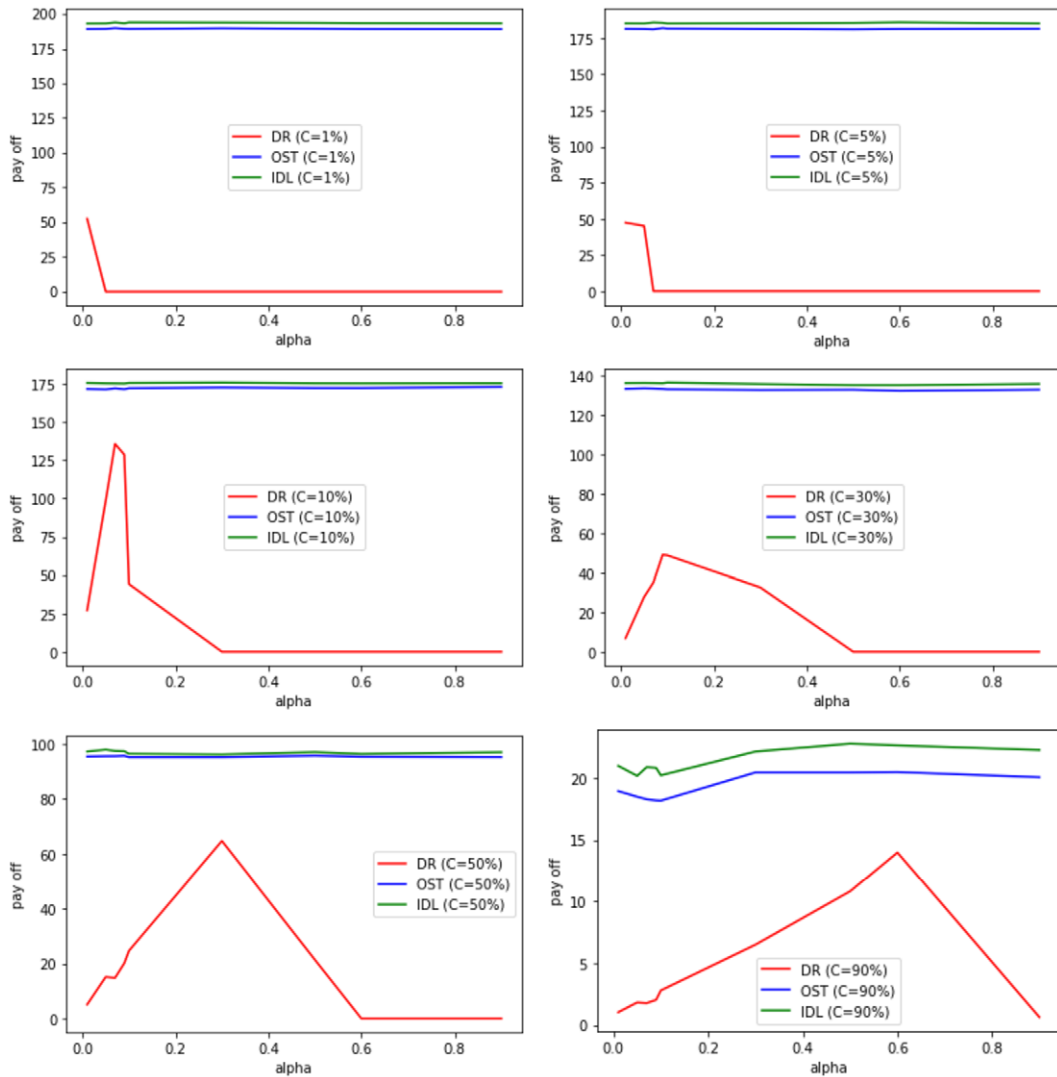


Fig. 5. (Experiment 2) The expected payoff vs delay cost with medium service demand rate  $\lambda = 10$  and tolerance  $\theta = 200$  for OST, DR and IDL models.

#### 4.2.3. Experiment 3

The settings for the Experiment 3 are:  $\lambda = 30$  and  $\theta = 600$  (high service demand rate). We performed the experiment with the defined  $\lambda$  and  $\theta$ . Our experiment is mainly affected by the cost  $C$  and  $\alpha$ . The experiment measures the payoff and the time that decision has been made. Moreover, the experiment is performed over three models. The first model is IDL which is being chosen to be the baseline to be compared with the other two models. The second model is our fundamental approach which is OST. Finally, the third model is DR. Figure 6 shows the different effects of the cost all over the three models and obviously shows the result of payoff and how our approach is as close as IDL compared with DR.

#### 4.3. Comparative assessment

The settings for the comparative assessment experiment 1 are:  $\lambda = 3$ ,  $\theta \in \{50, 100, 150, 200\}$  and  $C = (1\%, 10\%, 30\%) \cdot \lambda$  (low service demand rate). We performed the experiment with the defined  $\lambda$  and  $\alpha = 0.1$ . Our experiment is mainly affected by different values of  $\theta$  and  $C$ . Furthermore, the experiment measures the payoff

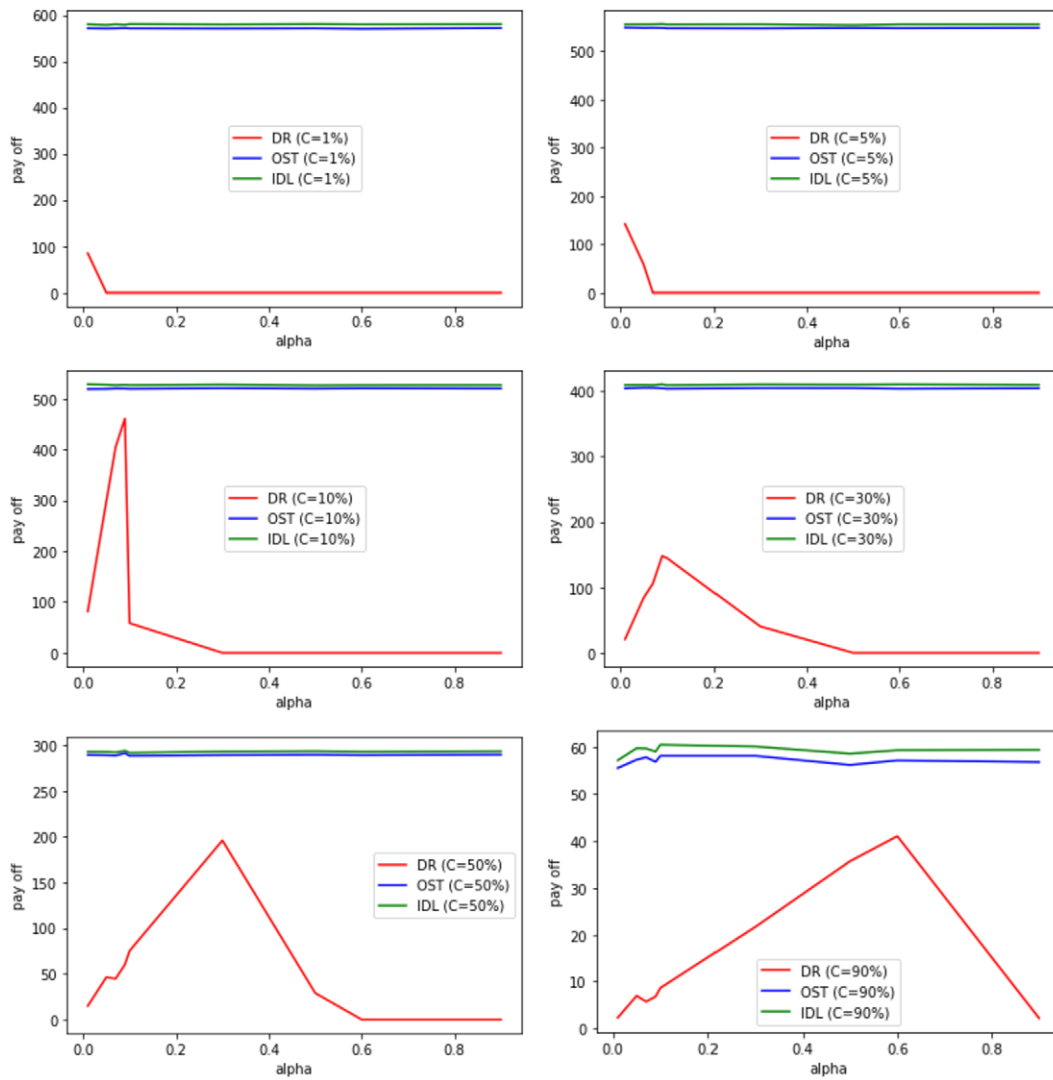


Fig. 6. (Experiment 3) The expected payoff vs delay cost with high service demand rate  $\lambda = 30$  and tolerance  $\theta = 600$  for OST, DR and IDL models.

$V^*$  and the time  $t^*$  that decision has been made. Moreover, the experiment was performed over four models. The first model is IDL which is being chosen to be the baseline. It used to be compared with the other three models. The second model is our fundamental approach which is OST. Then, the third model is DR. Finally, the fourth model is SECPRO. The experiment tested four different models. Figure 7 shows the different effects of  $\theta$  over all four models. Moreover, Fig. 7 obviously shows the result of payoff and how our approach is as close as IDL compared with DR and SECPRO.

The settings for the comparative assessment experiment 2 are:  $\lambda = 30$ ,  $\theta \in \{600, 1200, 1800, 2400\}$  and  $C = (1\%, 10\%, 30\%) \cdot \lambda$  (high service demand rate): Here, we performed the experiment with the defined  $\lambda$  and  $\alpha = 0.1$ . Our experiment is mainly affected by different values of  $\theta$  and  $C$ . Furthermore, the experiment measures the payoff  $V^*$  and the time  $t^*$  that decision has been made. Moreover, the experiment was performed over four models. The first model is IDL which is being chosen to be the baseline. It used to be compared with the other three models. The second model is our fundamental approach which is OST. Then, the third model is DR. Finally, the fourth model is SECPRO. The experiment tested four different models. Figure 8 shows the different effects of  $\theta$  over all four

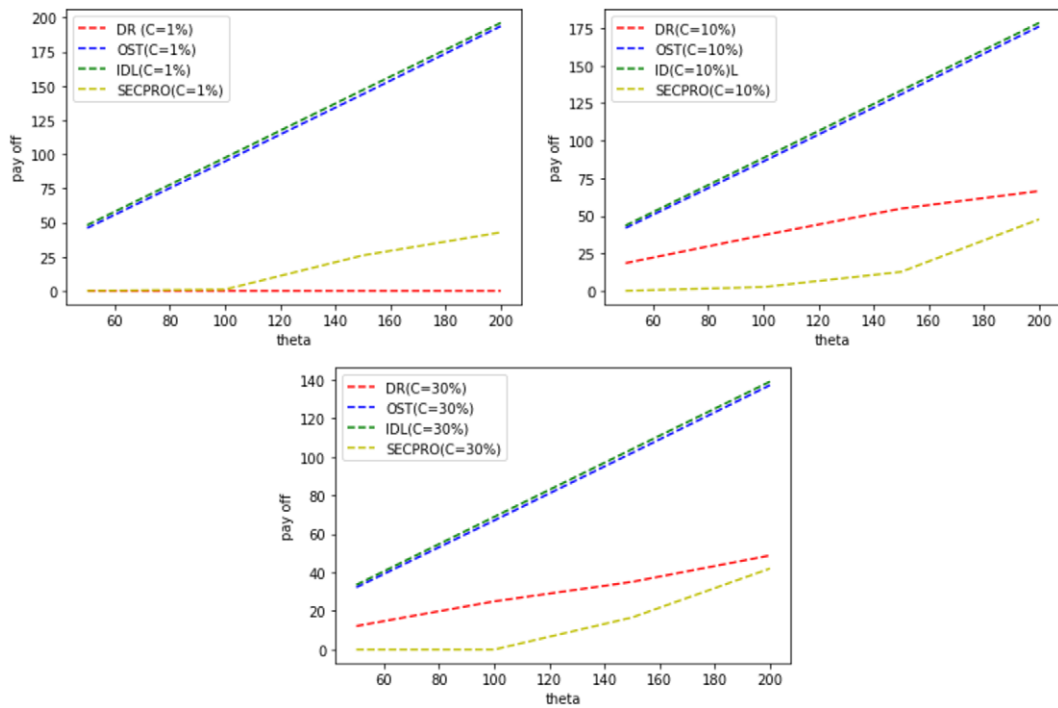


Fig. 7. Expected payoff vs tolerance  $\theta$  for diverse delay cot values with low service demand rate  $\lambda = 3$  and tolerance  $\theta \in \{50, 100, 150, 200\}$  for DR, OST, IDL and SECPRO models.

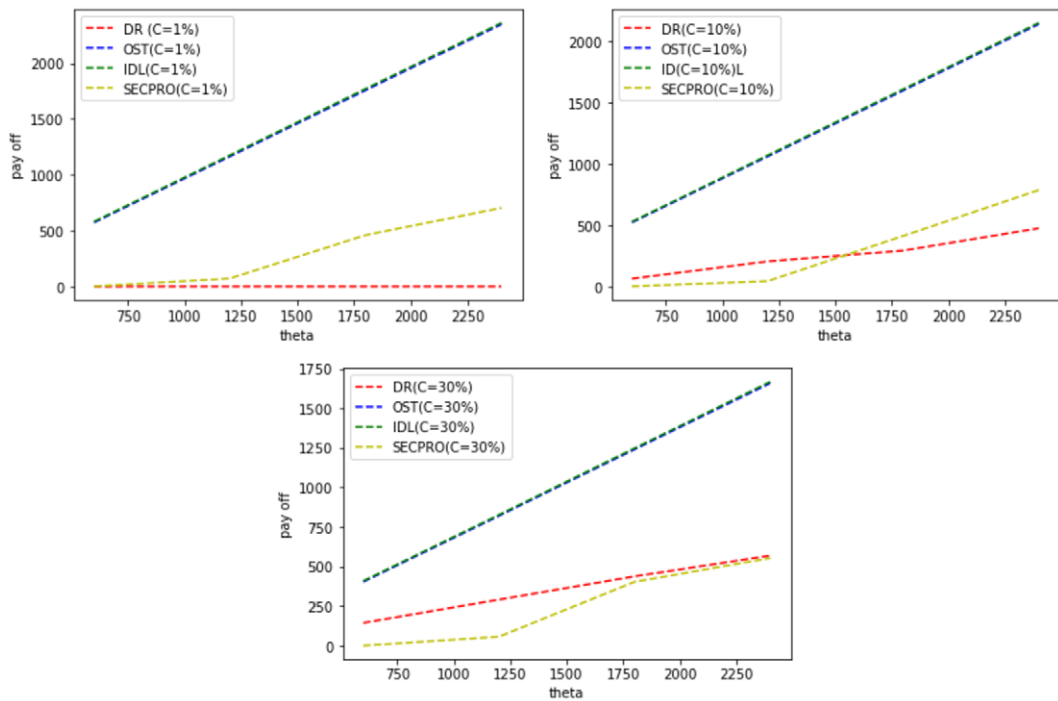


Fig. 8. Expected payoff vs tolerance  $\theta$  for diverse delay cost values with high service demand rate  $\lambda = 30$  and tolerance  $\theta \in \{600, 1200, 1800, 2400\}$  for DR, OST, IDL and SECPRO models.

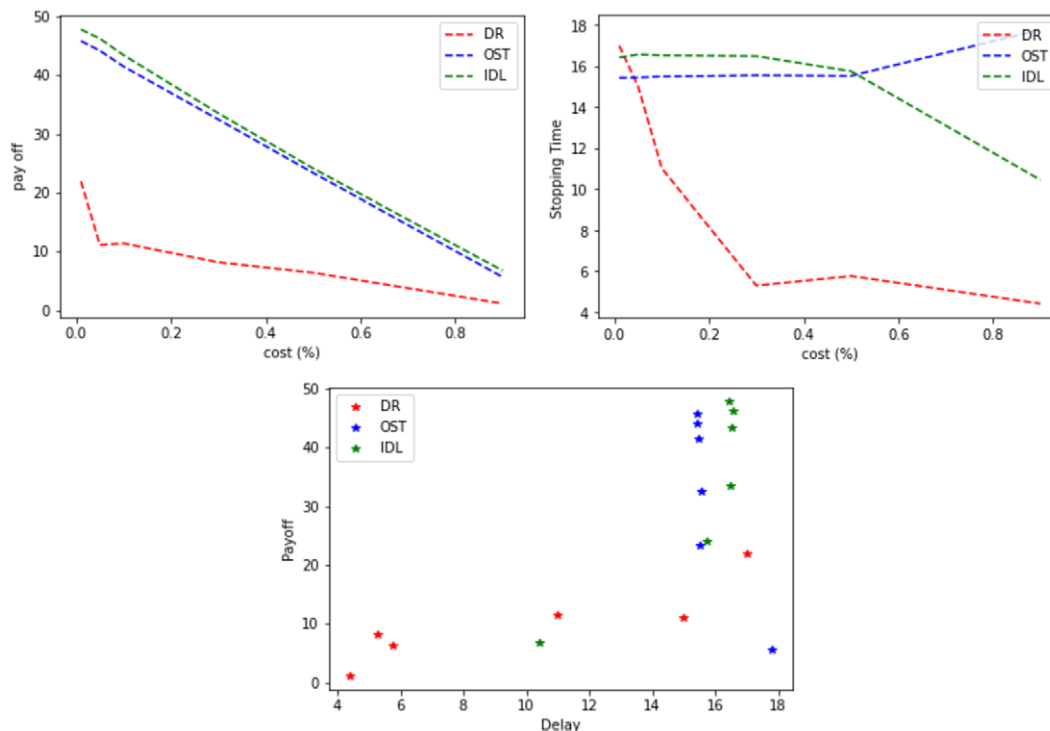


Fig. 9. Expected payoff vs cost and delay for OST, IDL (low service demand rate  $\lambda = 3$  and tolerance  $\theta = 50$ ).

models. Moreover, Fig. 8 obviously shows the result of payoff and how our approach is as close as IDL compared with DR and SECPRO.

In order to test the effectiveness of our theory (OST) in the first three experiments, where its results were compared with the IDL and another rule is DR. They were all tested on different parameters, which start from normal requests to medium requests and end with dense requests. Moreover, the results of our approach (OST) are amazing, which prove the effectiveness of the theory. As shown in Figs 9, 10 and 11, it obvious that Fig. 9 represents  $\lambda = 3$  and  $\theta = 50$ , as well as Fig. 10 shows  $\lambda = 10$  and  $\theta = 200$ . Also, Fig. 11 shows  $\lambda = 30$  and  $\theta = 600$ . In fact, each Figure contains three graphs, which represent the cost versus payoff, the cost versus stopping time, and finally delay versus payoff.

In Fig. 9, we see the results of our theory (OST), which are very close to IDL, and this is illustrated in the graphic that presents cost versus payoff. However, after increasing the cost to more than (0.5) cost in the figure that shows the cost versus the Stopping time, we find there is a small difference, and despite that, our approach obtains a delay time less than IDL in all results where the cost is (0.5) or less, and this is what makes it unique.

In Fig. 10, the receiving of requests were increased by making  $\lambda = 10$ , where we notice in the graph that shown the cost versus the payoff is an almost perfect match between our theory (OST) and IDL. Also, the delay time difference decreased after more than (0.5) the cost, however, we obtained a perfect stopping time between our theory (OST) and IDL as shown in the graph that illustrates the cost versus the stopping time as well as the graph showing the delay time versus the payoff.

Figure 11 shows the robustness and effectiveness of our theory (OST). It is shown clearly OST's ability to deal with heavy demands where  $\lambda = 30$ . The results were amazing even after increasing the cost by more than (0.5) and this is obvious in the graph that shows the cost versus the stopping time and the graph that shows the delay time against the payoff. As the delay in our theory (OST) is very close to the delay for IDL even after (0.5) cost. We also note in the graph that presents cost versus payoff, a great match between our theory (OST) and IDL.

Moreover, we tested the effectiveness of our theory (OST) in the fourth & fifth experiments, where its results were compared with the three methodologies which are IDL, DR, and SECPRO. They were all tested on two different

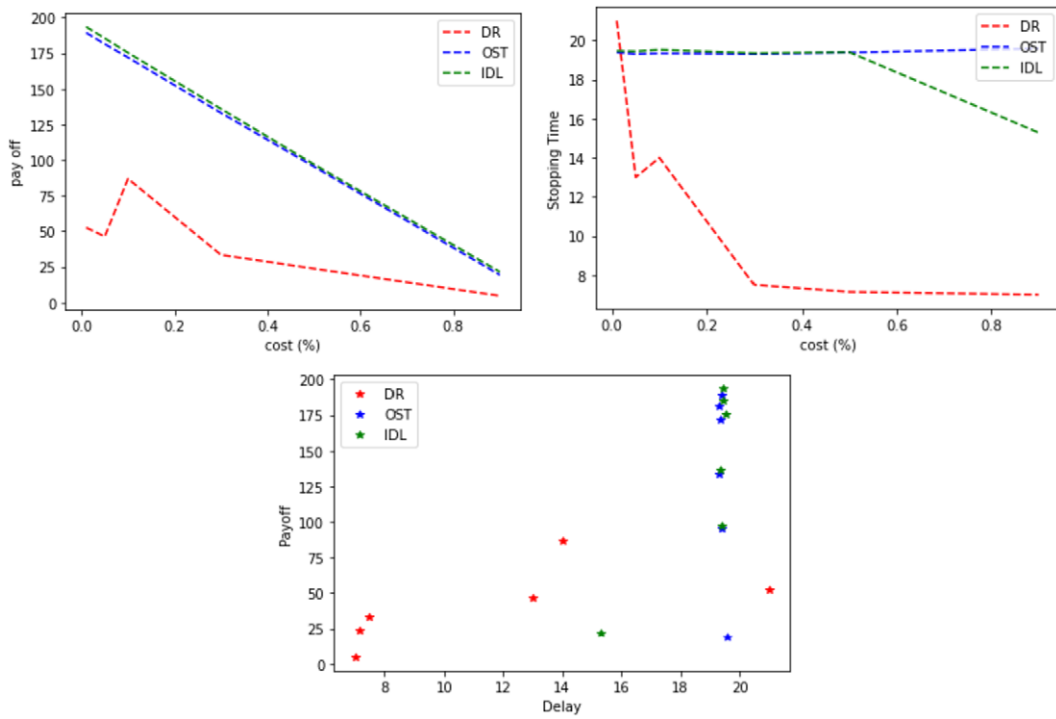


Fig. 10. Expected payoff vs cost and delay for OST, IDL (medium service demand rate  $\lambda = 10$  and tolerance  $\theta = 200$ ).

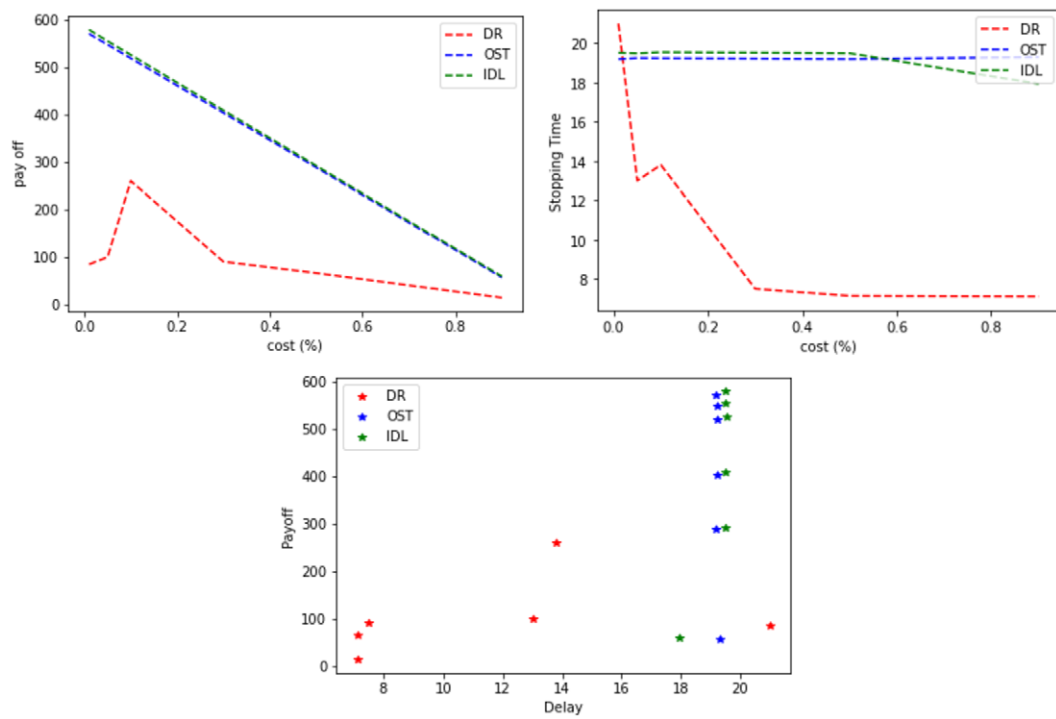


Fig. 11. Expected payoff vs cost and delay for OST, IDL (high service demand rate  $\lambda = 30$  and tolerance  $\theta = 600$ ).

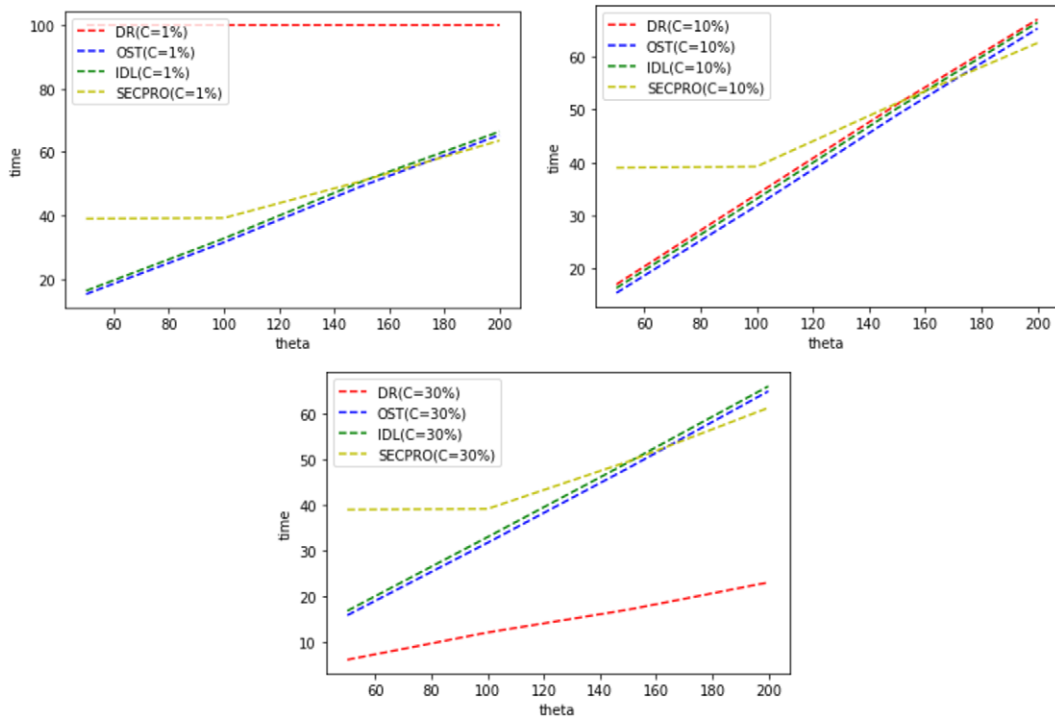


Fig. 12. Expected delay (time to decision) vs tolerance  $\theta$  for diverse cost values (low service demand rate  $\lambda = 3$  and tolerance  $\theta \in \{50, 100, 150, 200\}$ ).

parameters, which are low & high requests. Moreover, the results of our approach (OST) are brilliant, which proves the effectiveness of our theory OST. As shown in Figs 12 and 13, it is obvious that Fig. 12 represents  $\lambda = 3$ ,  $\theta \in \{50, 100, 150, 200\}$  and  $C = (1\%, 10\%, 30\%) \cdot \lambda$ , as well as Fig. 13 shows  $\lambda = 30$ ,  $\theta \in \{600, 1200, 1800, 2400\}$  and  $C = (1\%, 10\%, 30\%) \cdot \lambda$ . In fact, each figure contains three graphs, which represent the time of the decision being made  $t^*$  versus payoff.

## 5. Conclusions

We provide a superior decision-making approach for users in EC that is based on the Optimal Stopping Theory. When using and adopting the OST in time-optimized decision-making for a service replication (pull-action) in EC environments, we offer a comprehensive review of the process. Based on the results of our experimental assessments, the OST model performs more effectively than the other approaches, which are suitable for use in Edge nodes, and do not demand for a significant amount of resources. Moreover, the highest payoff is obtained by our models, which outperform alternative baseline solutions. In further work, our aim is to analyze how the characteristics of mobile applications and data timeliness affect our OST models and build new ones that take these limitations into account.

## Acknowledgements

The authors would like to thank the Saudi Prime Minister of Defense (Crown Prince Mohammad Bin Salman), the Royal Saudi Air Force, Senior Engineer Mohammad Aleissa, Saudi Arabia, and the Saudi Arabian Cultural Bureau in the UK for their support and encouragement.



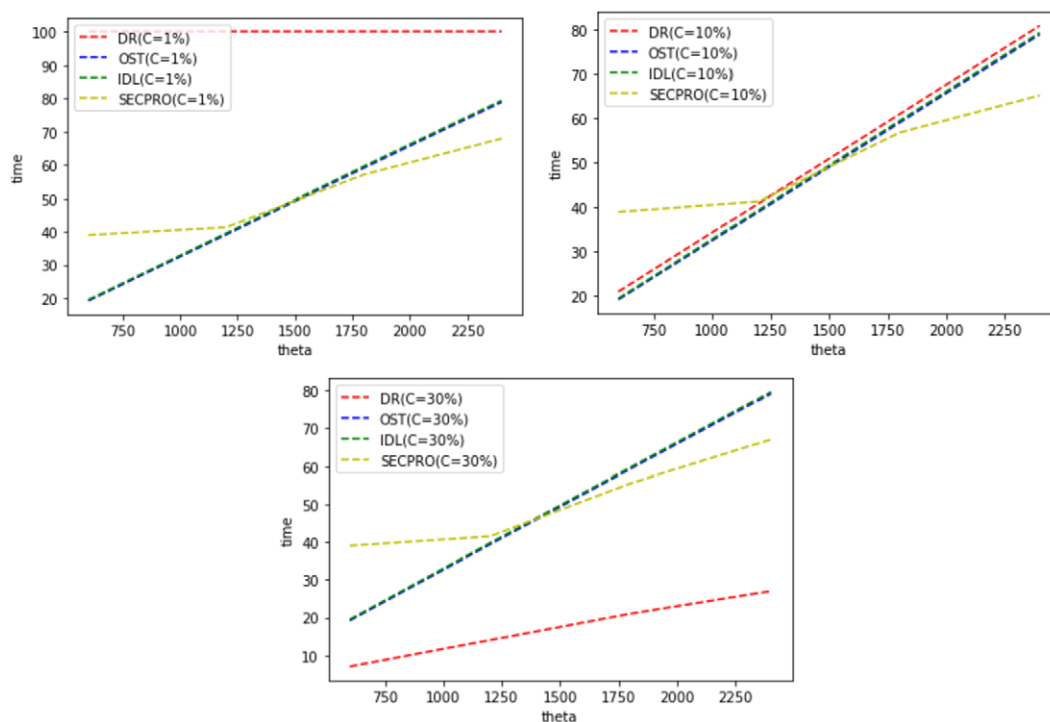


Fig. 13. Expected delay (time to decision) vs tolerance  $\theta$  for diverse cost values (high service demand rate  $\lambda = 30$  and tolerance  $\theta \in \{600, 1200, 1800, 2400\}$ ).

### Conflict of interest

The authors have no conflict of interest to report.

### References

- [1] M.G.R. Alam, M.M. Hassan, M.Z. Uddin, A. Almogren and G. Fortino, Autonomic computation offloading in mobile edge for iot applications, *Future Generation Computer Systems* **90** (2019), 149–157. doi:10.1016/j.future.2018.07.050.
- [2] M.G.R. Alam, Y.K. Tun and C.S. Hong, Multi-agent and reinforcement learning based code offloading in mobile fog, in: *2016 International Conference on Information Networking (ICOIN)*, IEEE, 2016, pp. 285–290.
- [3] T. Alfakih, M.M. Hassan, A. Gumaei, C. Savaglio and G. Fortino, Task offloading and resource allocation for mobile edge computing by deep reinforcement learning based on sarsa, *IEEE Access* **8** (2020), 54074–54084. doi:10.1109/ACCESS.2020.2981434.
- [4] I. Alghamdi, C. Anagnostopoulos and D.P. Pezaros, On the optimality of task offloading in mobile edge computing environments, in: *2019 IEEE Global Communications Conference (GLOBECOM)*, IEEE, 2019, pp. 1–6.
- [5] I. Alghamdi, C. Anagnostopoulos and D.P. Pezaros, Data quality-aware task offloading in mobile edge computing: An optimal stopping theory approach, *Future Generation Computer Systems* **117** (2021), 462–479. doi:10.1016/j.future.2020.12.017.
- [6] C. Anagnostopoulos, S. Hadjiefthymiades and K. Kolomvatsos, Accurate, dynamic, and distributed localization of phenomena for mobile sensor networks, *ACM Transactions on Sensor Networks (TOSN)* **12**(2) (2016), 1–59. doi:10.1145/2882966.
- [7] C. Anagnostopoulos and K. Kolomvatsos, A delay-resilient and quality-aware mechanism over incomplete contextual data streams, *Information Sciences* **355** (2016), 90–109. doi:10.1016/j.ins.2016.03.020.
- [8] P. Bellavista, A. Zanni and M. Solimando, A migration-enhanced edge computing support for mobile devices in hostile environments, in: *2017 13th International Wireless Communications and Mobile Computing Conference (IWCMC)*, IEEE, 2017, pp. 957–962. doi:10.1109/IWCMC.2017.7986415.
- [9] D. Callegaro and M. Levorato, Optimal edge computing for infrastructure-assisted uav systems, *IEEE Transactions on Vehicular Technology* **70**(2) (2021), 1782–1792. doi:10.1109/TVT.2021.3051378.
- [10] W. Chang, Y. Xiao, W. Lou and G. Shou, Offloading decision in edge computing for continuous applications under uncertainty, *IEEE Transactions on Wireless Communications* **19**(9) (2020), 6196–6209. doi:10.1109/TWC.2020.3001012.

- [11] D. Chemodanov, F. Esposito, A. Sukhov, P. Calyam, H. Trinh and Z. Oraibi, Agra: Ai-augmented geographic routing approach for iot-based incident-supporting applications, *Future Generation Computer Systems* **92** (2019), 1051–1065. doi:10.1016/j.future.2017.08.009.
- [12] M. Chen and Y. Hao, Task offloading for mobile edge computing in software defined ultra-dense network, *IEEE Journal on Selected Areas in Communications* **36**(3) (2018), 587–597. doi:10.1109/JSAC.2018.2815360.
- [13] Q. Chen, W. Wang, F. Wu, S. De, R. Wang, B. Zhang and X. Huang, A survey on an emerging area: Deep learning for smart city data, *IEEE Transactions on Emerging Topics in Computational Intelligence* **3**(5) (2019), 392–410. doi:10.1109/TETCI.2019.2907718.
- [14] S. Chen, H. Chen, J. Ruan and Z. Wang, Context-aware online offloading strategy with mobility prediction for mobile edge computing, in: *2021 International Conference on Computer Communications and Networks (ICCCN)*, IEEE, 2021, pp. 1–9.
- [15] Y. Chen, N. Zhang, Y. Zhang, X. Chen, W. Wu and X. Shen, Energy efficient dynamic offloading in mobile edge computing for Internet of things, *IEEE Transactions on Cloud Computing* **9**(3) (2019), 1050–1060. doi:10.1109/TCC.2019.2898657.
- [16] T.Q. Dinh, Q.D. La, T.Q. Quek and H. Shin, Learning for computation offloading in mobile edge computing, *IEEE Transactions on Communications* **66**(12) (2018), 6353–6367. doi:10.1109/TCOMM.2018.2866572.
- [17] T.Q. Dinh, J. Tang, Q.D. La and T.Q. Quek, Offloading in mobile edge computing: Task allocation and computational frequency scaling, *IEEE Transactions on Communications* **65**(8) (2017), 3571–3584.
- [18] O. Elijah, T.A. Rahman, I. Orikumhi, C.Y. Leow and M.N. Hindia, An overview of Internet of things (iot) and data analytics in agriculture: Benefits and challenges, *IEEE Internet of things Journal* **5**(5) (2018), 3758–3773. doi:10.1109/JIOT.2018.2844296.
- [19] A.E. Eshratifar and M. Pedram, Energy and performance efficient computation offloading for deep neural networks in a mobile cloud computing environment, in: *Proceedings of the 2018 on Great Lakes Symposium on VLSI*, 2018, pp. 111–116. doi:10.1145/3194554.3194565.
- [20] Z. Gao, Q. Jiao, K. Xiao, Q. Wang, Z. Mo and Y. Yang, Deep reinforcement learning based service migration strategy for edge computing, in: *2019 IEEE International Conference on Service-Oriented System Engineering (SOSE)*, IEEE, 2019, pp. 116–1165. doi:10.1109/SOSE.2019.00025.
- [21] Q. Gu, Y. Jian, G. Wang, R. Fan, H. Jiang and Z. Zhong, Mobile edge computing via wireless power transfer over multiple fading blocks: An optimal stopping approach, *IEEE Transactions on Vehicular Technology* **69**(9) (2020), 10348–10361. doi:10.1109/TVT.2020.3005406.
- [22] N. Hassan, S. Gillani, E. Ahmed, I. Yaqoob and M. Imran, The role of edge computing in Internet of things, *IEEE communications magazine* **56**(11) (2018), 110–115. doi:10.1109/MCOM.2018.1700906.
- [23] X. Huang, K. Xu, C. Lai, Q. Chen and J. Zhang, Energy-efficient offloading decision-making for mobile edge computing in vehicular networks, *EURASIP Journal on Wireless Communications and Networking* **2020**(1) (2020), 1–16. doi:10.1186/s13638-019-1618-7.
- [24] N. Kaur, A. Mittal, A. Kumar and R. Kumar, Healthcare monitoring through fog computing: A survey, *ECS Transactions* **107**(1) (2022), 7689. doi:10.1149/10701.7689ecst.
- [25] M. Ke, Z. Gao, Y. Wu, X. Gao and K.-K. Wong, Massive access in cell-free massive mimo-based Internet of things: Cloud computing and edge computing paradigms, *IEEE Journal on Selected Areas in Communications* **39**(3) (2020), 756–772. doi:10.1109/JSAC.2020.3018807.
- [26] K. Kolomvatsos and C. Anagnostopoulos, A proactive statistical model supporting services and tasks management in pervasive applications, *IEEE Transactions on Network and Service Management* (2022).
- [27] K. Kolomvatsos, C. Anagnostopoulos and S. Hadjiefthymiades, An efficient environmental monitoring system adopting data fusion, prediction, & fuzzy logic, in: *2015 6th International Conference on Information, Intelligence, Systems and Applications (IISA)*, IEEE, 2015, pp. 1–6.
- [28] K. Kolomvatsos, C. Anagnostopoulos and S. Hadjiefthymiades, Data fusion and type-2 fuzzy inference in contextual data stream monitoring, *IEEE Transactions on Systems, Man, and Cybernetics: Systems* **47**(8) (2016), 1839–1853. doi:10.1109/TSMC.2016.2560533.
- [29] K. Kolomvatsos, C. Anagnostopoulos and S. Hadjiefthymiades, Distributed localized contextual event reasoning under uncertainty, *IEEE Internet of Things Journal* **4**(1) (2016), 183–191. doi:10.1109/JIOT.2016.2638119.
- [30] K. Kolomvatsos, C. Anagnostopoulos, M. Koziri and T. Loukopoulos, Proactive & time-optimized data synopsis management at the edge, *IEEE Transactions on Knowledge and Data Engineering* (2020).
- [31] L. Kuang, T. Gong, S. OuYang, H. Gao and S. Deng, Offloading decision methods for multiple users with structured tasks in edge computing for smart cities, *Future Generation Computer Systems* **105** (2020), 717–729. doi:10.1016/j.future.2019.12.047.
- [32] Z. Kuang, L. Li, J. Gao, L. Zhao and A. Liu, Partial offloading scheduling and power allocation for mobile edge computing systems, *IEEE Internet of Things Journal* **6**(4) (2019), 6774–6785. doi:10.1109/JIOT.2019.2911455.
- [33] Z. Liang, Y. Liu, T.-M. Lok and K. Huang, Service migration for multi-cell mobile edge computing, in: *GLOBECOM 2020-2020 IEEE Global Communications Conference*, IEEE, 2020, pp. 1–6.
- [34] C.-C. Liao, T.-S. Chen and A.-Y. Wu, Real-time multi-user detection engine design for iot applications via modified sparsity adaptive matching pursuit, *IEEE Transactions on Circuits and Systems I: Regular Papers* **66**(8) (2019), 2987–3000. doi:10.1109/TCSI.2019.2903193.
- [35] F. Messaoudi, A. Ksentini and P. Bertin, On using edge computing for computation offloading in mobile network, in: *GLOBECOM 2017-2017 IEEE Global Communications Conference*, IEEE, 2017, pp. 1–7.
- [36] T. Ouyang, R. Li, X. Chen, Z. Zhou and X. Tang, Adaptive user-managed service placement for mobile edge computing: An online learning approach, in: *IEEE INFOCOM 2019-IEEE Conference on Computer Communications*, IEEE, 2019, pp. 1468–1476. doi:10.1109/INFOCOM.2019.8737560.
- [37] T. Ouyang, Z. Zhou and X. Chen, Follow me at the edge: Mobility-aware dynamic service placement for mobile edge computing, *IEEE Journal on Selected Areas in Communications* **36**(10) (2018), 2333–2345. doi:10.1109/JSAC.2018.2869954.
- [38] Y. Pan, M. Chen, Z. Yang, N. Huang and M. Shikh-Bahaei, Energy-efficient noma-based mobile edge computing offloading, *IEEE Communications Letters* **23**(2) (2018), 310–313. doi:10.1109/LCOMM.2018.2882846.
- [39] G. Peskir and A. Shiryaev, *Optimal Stopping and Free-Boundary Problems*, Springer, 2006.

- [40] Q. Qi and F. Tao, A smart manufacturing service system based on edge computing, fog computing, and cloud computing, *IEEE Access* **7** (2019), 86769–86777. doi:[10.1109/ACCESS.2019.2923610](https://doi.org/10.1109/ACCESS.2019.2923610).
- [41] Q. Qi, J. Wang, Z. Ma, H. Sun, Y. Cao, L. Zhang and J. Liao, Knowledge-driven service offloading decision for vehicular edge computing: A deep reinforcement learning approach, *IEEE Transactions on Vehicular Technology* **68**(5) (2019), 4192–4203. doi:[10.1109/TVT.2019.2894437](https://doi.org/10.1109/TVT.2019.2894437).
- [42] M. Sato, Y. Fukuyama, T. Iizaka and T. Matsui, Total optimization of energy networks in a smart city by multi-swarm differential evolutionary particle swarm optimization, *IEEE Transactions on Sustainable Energy* **10**(4) (2018), 2186–2200. doi:[10.1109/TSTE.2018.2882203](https://doi.org/10.1109/TSTE.2018.2882203).
- [43] S. Singh, I.-H. Ra, W. Meng, M. Kaur and G.H. Cho, Sh-blockcc: A secure and efficient Internet of things smart home architecture based on cloud computing and blockchain technology, *International Journal of Distributed Sensor Networks* **15**(4) (2019). doi:[10.1177/1550147719844159](https://doi.org/10.1177/1550147719844159).
- [44] F. Sun, F. Hou, N. Cheng, M. Wang, H. Zhou, L. Gui and X. Shen, Cooperative task scheduling for computation offloading in vehicular cloud, *IEEE Transactions on Vehicular Technology* **67**(11) (2018), 11049–11061. doi:[10.1109/TVT.2018.2868013](https://doi.org/10.1109/TVT.2018.2868013).
- [45] M. Sun, Y. Wang, G. Strbac and C. Kang, Probabilistic peak load estimation in smart cities using smart meter data, *IEEE Transactions on Industrial Electronics* **66**(2) (2018), 1608–1618. doi:[10.1109/TIE.2018.2803732](https://doi.org/10.1109/TIE.2018.2803732).
- [46] Y. Sun, X. Guo, S. Zhou, Z. Jiang, X. Liu and Z. Niu, Learning-based task offloading for vehicular cloud computing systems, in: *2018 IEEE International Conference on Communications (ICC)*, IEEE, 2018, pp. 1–7.
- [47] Y. Sun, S. Zhou and J. Xu, Emm: Energy-aware mobility management for mobile edge computing in ultra dense networks, *IEEE Journal on Selected Areas in Communications* **35**(11) (2017), 2637–2646. doi:[10.1109/JSAC.2017.2760160](https://doi.org/10.1109/JSAC.2017.2760160).
- [48] L. Tang and S. He, Multi-user computation offloading in mobile edge computing: A behavioral perspective, *IEEE Network* **32**(1) (2018), 48–53. doi:[10.1109/MNET.2018.1700119](https://doi.org/10.1109/MNET.2018.1700119).
- [49] M. Tang and V.W.S. Wong, Deep reinforcement learning for task offloading in mobile edge computing systems, *IEEE Transactions on Mobile Computing* (2020).
- [50] F. Wang, J. Xu, X. Wang and S. Cui, Joint offloading and computing optimization in wireless powered mobile-edge computing systems, *IEEE Transactions on Wireless Communications* **17**(3) (2017), 1784–1797. doi:[10.1109/TWC.2017.2785305](https://doi.org/10.1109/TWC.2017.2785305).
- [51] J. Wang, J. Pan, F. Esposito, P. Calyam, Z. Yang and P. Mohapatra, Edge cloud offloading algorithms: Issues, methods, and perspectives, *ACM Computing Surveys (CSUR)* **52**(1) (2019), 1–23. doi:[10.1145/3214306](https://doi.org/10.1145/3214306).
- [52] S. Wang, Y. Guo, N. Zhang, P. Yang, A. Zhou and X. Shen, Delay-aware microservice coordination in mobile edge computing: A reinforcement learning approach, *IEEE Transactions on Mobile Computing* **20**(3) (2019), 939–951. doi:[10.1109/TMC.2019.2957804](https://doi.org/10.1109/TMC.2019.2957804).
- [53] S. Wang, R. Urgaonkar, M. Zafer, T. He, K. Chan and K.K. Leung, Dynamic service migration in mobile edge computing based on Markov decision process, *IEEE/ACM Transactions on Networking* **27**(3) (2019), 1272–1288. doi:[10.1109/TNET.2019.2916577](https://doi.org/10.1109/TNET.2019.2916577).
- [54] S. Wang, J. Xu, N. Zhang and Y. Liu, A survey on service migration in mobile edge computing, *IEEE Access* **6** (2018), 23511–23528. doi:[10.1109/ACCESS.2018.2828102](https://doi.org/10.1109/ACCESS.2018.2828102).
- [55] Y. Wang, M. Sheng, X. Wang, L. Wang and J. Li, Mobile-edge computing: Partial computation offloading using dynamic voltage scaling, *IEEE Transactions on Communications* **64**(10) (2016), 4268–4282. doi:[10.1109/TCOMM.2016.2594794](https://doi.org/10.1109/TCOMM.2016.2594794).
- [56] J. Xu, X. Ma, A. Zhou, Q. Duan and S. Wang, Path selection for seamless service migration in vehicular edge computing, *IEEE Internet of Things Journal* **7**(9) (2020), 9040–9049. doi:[10.1109/JIOT.2020.3000300](https://doi.org/10.1109/JIOT.2020.3000300).
- [57] X. Xu, Q. Huang, X. Yin, M. Abbasi, M.R. Khosravi and L. Qi, Intelligent offloading for collaborative smart city services in edge computing, *IEEE Internet of Things Journal* **7**(9) (2020), 7919–7927. doi:[10.1109/JIOT.2020.3000871](https://doi.org/10.1109/JIOT.2020.3000871).
- [58] X. Xu, X. Liu, Z. Xu, F. Dai, X. Zhang and L. Qi, Trust-oriented iot service placement for smart cities in edge computing, *IEEE Internet of Things Journal* **7**(5) (2019), 4084–4091. doi:[10.1109/JIOT.2019.2959124](https://doi.org/10.1109/JIOT.2019.2959124).
- [59] Z. Yu, Y. Gong, S. Gong and Y. Guo, Joint task offloading and resource allocation in uav-enabled mobile edge computing, *IEEE Internet of Things Journal* **7**(4) (2020), 3147–3159. doi:[10.1109/JIOT.2020.2965898](https://doi.org/10.1109/JIOT.2020.2965898).
- [60] Q. Yuan, J. Li, H. Zhou, T. Lin, G. Luo and X. Shen, A joint service migration and mobility optimization approach for vehicular edge computing, *IEEE Transactions on Vehicular Technology* **69**(8) (2020), 9041–9052. doi:[10.1109/TVT.2020.2999617](https://doi.org/10.1109/TVT.2020.2999617).
- [61] H. Zhang, Z. Wang and K. Liu, V2x offloading and resource allocation in sdn-assisted mec-based vehicular networks, *China Communications* **17**(5) (2020), 266–283. doi:[10.23919/JCC.2020.05.020](https://doi.org/10.23919/JCC.2020.05.020).
- [62] J. Zhang, Y. Wang, S. Li and S. Shi, An architecture for iot-enabled smart transportation security system: A geospatial approach, *IEEE Internet of Things Journal* **8**(8) (2020), 6205–6213. doi:[10.1109/JIOT.2020.3041386](https://doi.org/10.1109/JIOT.2020.3041386).
- [63] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo and J. Zhang, Edge intelligence: Paving the last mile of artificial intelligence with edge computing, *Proceedings of the IEEE* **107**(8) (2019), 1738–1762. doi:[10.1109/JPROC.2019.2918951](https://doi.org/10.1109/JPROC.2019.2918951).