

A cache-based approach toward improved scheduling in fog computing

Osama Amir Khan¹ | Saif U. R. Malik²  | Faizan M. Baig¹ | Saif Ul Islam³ | Haris Pervaiz⁴  | Hassan Malik⁵ | Syed Hassan Ahmed⁶

¹Department of Computer Science, COMSATS University, Islamabad, Pakistan

²Cybernetica AS, Tallinn, Estonia

³Dr. A. Q. Khan Institute of Computer Science and Information Technology, Rawalpindi, Pakistan

⁴School of Computing and Communications (SCC), Lancaster University, Lancaster, UK

⁵Thomas Johann Seebeck Department of Electronics, Tallinn University of Technology, Tallinn, Estonia

⁶Department of Computer Science, Georgia Southern University, Statesboro, Georgia

Correspondence

Saif U. R. Malik, Cybernetica AS
Tallinn, Estonia.
Email: saif.rehmanmalik@cyber.ee

Abstract

Fog computing is a promising technique to reduce the latency and power consumption issues of the Internet of Things (IoT) ecosystem by enabling storage and computational resource close to the end-user devices with additional benefits such as improved execution time and processing. However, with an increase in IoT devices, the resource allocation and job scheduling became a complicated and cumbersome task due to limited and heterogeneous resources along with the locality restriction in such computing environment. Therefore, this paper proposes a cache-based approach for efficient resource allocation in fog computing environment, while maintaining the quality of service. The proposed algorithm is realized using iFogSim simulator and a comprehensive comparison is presented with the traditional First Come First Served and Shortest Job First policies. The performance evaluation revealed that with the proposed scheme the execution time, latency, processing delays and power consumption decreased by 38%, 11.1%, 6%, and 17.8%, respectively, as compared to those of the traditional schemes.

KEYWORDS

cache, cloud computing, fog computing, IoT, job scheduling, QoE, QoS

1 | INTRODUCTION

Internet of Things (IoT) envisioned to enable an ecosystem where physical things embedded with sensors and communication technologies can exchange data among each other and with the network to enables services such as smart homes, smart city, smart health care, connected vehicles, etc.^{1,2} The basic IoT architecture constitutes of four components namely Things, Gateway, Communication Network, and Cloud Platform.² Things are commonly referred as IoT end devices enabled with sensor or actuator for data collection. The collected data from these devices are then forwarded to Gateways, nowadays, named as Fog Nodes. They act as an intermediate node between IoT devices and Cloud Platform providing the needed connectivity, scalability, security, and manageability. However, communication network refers to both wired and wireless communication technology needed for the actual transfer of information among these components. Lastly, the Cloud Platform can be regarded as main data collection center constitute of large number of servers and storage capabilities that are linked together. The key functionality of Cloud Platform is to process large amount of data using cloud computing techniques and provide meaningful information that can be used to support IoT application and services.

However, the rapid advancement in IoT ecosystem is expected to have billions of connected devices in future with more stringent application requirements in terms of latency, processing delay, power consumption and execution time, etc. To meet such requirements, new computing technologies such as edge or fog computing are introduced to partially or completely process the data at the device or gateway level.³ Fog computing inherits the reduced functionalities of cloud computing techniques that can be executed on resource constraint devices such as gateways. Fog computing can be viewed as geographically distributed computing paradigm, having heterogeneous devices at the edge of the network that are connected collaboratively to give elastic computation, communication, and storage services.⁴ Fog computing not only reduces the execution time and latency by providing fast data processing but also reduces the burden on communication network and power required by devices for long communication. This will allow the communication networks to serve more devices as well.

Despite having promising features, fog computing still faces high latency issues, as reported in Reference 5 due to the lack of efficient resource and job scheduling algorithm. There are several factors that makes a resource allocation in a fog computing environment a challenging task, such as resource scarcity, heterogeneity, geographic restrictions, and varying resource demands.⁶ The goals of scheduling and resource allocation is to increase the efficiency of the use of resources, satisfy the quality of service (QoS) requirements, meanwhile maximizing the profit of both fog nodes and user devices.⁷

Therefore, this paper presents a cache-based approach (CBA) for optimal resource and job scheduling in fog computing environment to achieve reduction in execution time, latency, internal processing delay, and power consumption. Caching is one of the promising technology for speeding up data retrieval time, reducing the number of path lengths, and improving system efficiency. Caching strategies exploit storage capacity to absorb traffic by replicating the most popular content closer to the node and enable low processing cost and remove a single point-of-failure.⁸⁻¹⁰ Therefore, the integration of caching in fog computing will enable fog nodes to identify the demand of the user and proactively select the most suitable contents to cache in geo-distributed nodes and improve the resource utilization. The main contribution of this paper can be summarized as follows:

- Firstly, the resource allocation in fog computing is modeled as delay minimization problem with constraints of available resources and number of jobs.
- Secondly, a CBA-based scheduling is proposed to address the formulated problem and to improve the efficiency of resource utilization while maintaining the QoS.
- Thirdly, a comprehensive comparison of the proposed scheme with the state-of-art schemes such as First Come First Serve (FCFS) and Shortest Job First (SJF) is presented. Moreover, the corresponding performance evaluation in terms of execution time, latency, internal processing delay, and power consumption are discussed in detail.

The remaining of the paper is organized as follows. Section 2 presents the related work on CBA for fog computing. Section 3 presents the fog computing architecture and its important components along with the interdependencies. The proposed scheduling methodology is presented in Section 4. Section 5 presents the performance evaluation of the proposed scheme and comprehensive comparison other traditional schemes in Section 5. Finally, the conclusion is drawn in Section 6.

2 | RELATED WORK

This section presents an overview of resource allocation schemes for resource constraints fog computing devices and highlight the limitation of the proposed schemes.

In the literature, one of the promising solutions for resource allocation for fog computing is by mean of smart gateway.¹¹ The concept is to use a smart gateway between an end device and cloud, that is capable to preprocess data, that is, filtering, monitoring, management, and resource allocation. Based on initial data processing, smart gateway can decide whether data need to be transmitted to cloud or the gateway itself can process the data. However, the resource optimization at smart gateway is still a challenging task. In this regard, in Reference 12, an author presents QoS-based resource allocation using Particle Swarm Optimization. The proposed algorithm improves the performance in terms of resource utilization and reduce delay. Moreover, in Reference 13, a bio-inspired algorithm based on Bees is proposed to optimize the task distribution among resource of end devices, fog nodes, and cloud server. The aim of the proposed

algorithm is to optimize central processing unit execution time and memory usage. Similarly, in Reference 14, load balancing algorithm between end device, fog node, and cloud is presented. The proposed algorithm aims to optimize execution time, resource allocation and deadlines based on heuristic approach. A threshold to control the number of jobs request is set for the fog layer. Once the threshold expires, the corresponding tasks forwarded to cloud layer for execution.

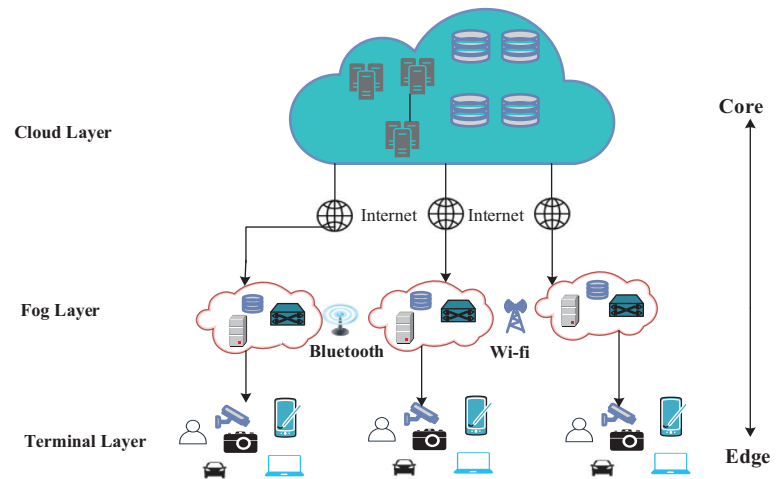
The architecture of fog computing to assist systems like 5G networks to attain high performance by ensuring optimal scheduling of job is presented in Reference 15. In this work, three policies are considered for job scheduling. In the first policy, fog node is randomly selected to execute the job from a uniform distribution called random policy. In second policy, fog nodes provide low latency which depends on systems current state and known as latency policy. The last policy is the capacity policy which selects fog node having most extreme number of outstanding resources among the candidate nodes. The simulation demonstrated that the least latency policy gives a superior outcome because of the accessibility of resources. The authors of the paper inferred that combination of the three policies together locate the most reasonable node for the job. Therefore, utilizing a solitary strategy may not be the best answer for the entire system. In Reference 16, authors present a load balancing mechanism in fog computing in which task distribution depends on graph partition. In this mechanism, tasks are allocated to multiple or single nodes of virtual machines depending on the requirements of the task. By using nondirectional graphs, physical nodes of the fog computing are represented in this work. These physical nodes came into a lot of virtual machine nodes as indicated by the accessible fog computing resources, where the virtual machine nodes give administrations to the users by means of graphics partition. To achieve this, the entire graph is used to create a minimum spanning tree and those edges that did not offer abundant resources are removed from the tree. The result of this graph shows the load balancing partition which is fingered by fog computing. Task runtime is achieved using this mechanism. However, the limitation of this technique is that for dynamic load balancing high performance is not achieved due to the regular repartitioning that is expected to deal with fog changes. Furthermore, in Reference 17, authors proposed mechanism for task scheduling and resource allocation that is based on container. The method is proposed to decrease the delay in execution of the task. Authors in Reference 18 proposed task scheduling algorithm in fog computing that is based on priority levels. The proposed algorithm consists of two steps. First step is assigning the task to nearest fog server. In second step all the requests are process in the three-priority queue within a fog server and reallocate the task to other fog servers if the selected fog server have insufficient resources. Finally, task is sent to cloud if fog layer does not have any resources.

Most of the above-mentioned studies focused to increase the performance by minimizing latency and execution time and ignored the resource re-allocation. This paper, combine the concept of smart gateway and caching to address the resource allocation problem for fog computing and presents CBA scheduling. The main idea of the proposed scheme is to incorporate cache module within smart gateway that stores the job and server information. The existence of the cache module will decrease the average waiting time of the jobs, which in return will have significant impact on the overall performance of the system. The presence of the smart gateway increases the performance, minimize latency and propagation delay. Moreover, the caches record will be used to reallocate the job to the cloud and fog nodes. The main objective is to assign the job to the most optimal resource that take minimum time and power to execute the job. This will improve the response time and reduce the cost in term of processing as well.

To the best of author's knowledge, the use of the caching for job scheduling in fog computing is under studied. In Reference 19, a community-based caching approach (CC) to solve cache pollution and cache monopoly problems in cloud computing-based high-performance web services. CC performance is compared with 13 other policies that are managed by cache, and in results its concluded that CC is better than other policies by achieving the cache-hit rate between 0.7% and 55%. The motivation of this research work is due to the appearance of caching as a cloud service, which supports web services with increasing user demands on its backend database servers. Similarly, in Reference 20 a Cost Aware Cache Replacement Policy for fog computing is proposed. The proposed algorithm aims to minimize the cache miss cost in a hybrid memory system.

3 | FOG COMPUTING ARCHITECTURE

Fog computing architecture consists of Cloud layer, Fog layer, and the Terminal layer are shown in Figure 1. Cloud layer includes multiple data storing devices and high-performance servers, and issue different application services like smart

FIGURE 1 Fog computing architecture

transportation, smart home, smart factory, etc. Cloud layer has huge storage and powerful computing capabilities to support extensive computing analysis and permanent storage of a huge amount of data. However, unlike traditional cloud computing architecture, not all computing and storage tasks go through the cloud. Depending on the demand load, the central modules of the cloud are managed and planned efficiently some control strategies to improve the use of cloud resources.²¹

Fog layer has numerous fog nodes, and each node includes a base station, routers, switches, gateway, access point, specific fog servers, etc. Fog node is a bridge between the cloud and end devices. Fog nodes can be static or mobile on a moving carrier. The end devices get services from fog nodes. Real-time analysis and low latency can be achieved in fog layer. Also fog layer is connected by IP core network with cloud data center. Fog nodes are connected with cloud to get more storage and computing capabilities. The third layer is a terminal layer and is closer to the end user and the physical environment. It includes different IoT gadgets like smart vehicles, mobile phones, sensors, smart card, and so on. These devices (sensors) are responsible for sensing the data of physical objects or events and transmit it to upper layer for processing and storage.²¹

The end user can directly communicate with fog server using wireless connections that are 4G LTE devices, Wi-Fi, LPWAN technologies and Bluetooth, etc. Cloud and fog server communicate to each other through a wired and wireless connection to access more application tools and computing services or resources. Mostly data stored in fog node is on brief premise. The cloud is more suitable for long-term data storage because of the availability of more resources than fog nodes. When the data is sent to the cloud, it is then not required to store on fog nodes.²²

The issue of latency between end user device and the cloud is addressed by fog computing architecture. Fog computing extends the cloud computing by moving the storage and computation resources at the edge of the network. Fog computing offers advantages like a fast response to delay-delicate applications, data aggregation for heterogeneous devices, gives data security and protection for sensitive data, avoids pointless communication by filtering the data before sending it to the cloud and provide context-aware and location-aware services. Along with these benefits, there are several challenges that still need to be addressed. These challenges include fog-cloud collaboration, service scalability (horizontal and vertical), fog scalability, fog resource management, and fog-based dedicated applications.²²

There are different performance metrics on which fog computing performance is measured. In Reference 22, researchers evaluated fog computing performance against different metrics of performance such as processing costs, processing delay, and processing power to show the gain in performance. Quality of experience (QoE) is another performance metrics to evaluate the performance of scheduling algorithm for fog computing.¹⁴

This paper proposes an extension of the fog computing architecture by introducing smart gateway enabled with caching for efficient resource allocation and evaluate the performance in terms of execution time, processing delay, latency, and power consumption. The detailed proposed architecture is shown in Figure 2 along with the detailed methodology in Section 4.

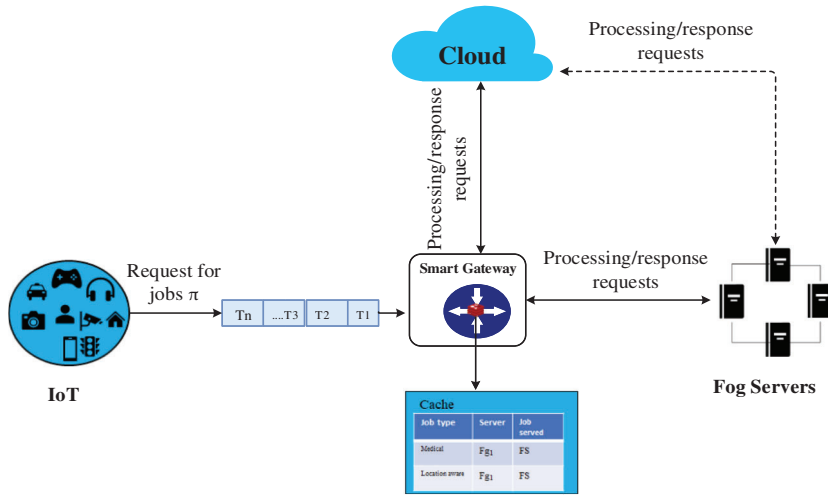


FIGURE 2 Cache-based approach

TABLE 1 List of Notations

T_{list}	List of job requested by all users	L	Length of job/tuple
T_i	Job requested by single user	R	Link bandwidth (bits/seconds)
C_{type}	Communication Type	D_c	Cloud data center
F_{ls}	Local storage of fog server	I_{pt}	Internal processing time
T_{exec}	Execution time	S_{type}	Service type requested by IoT
T_{type}	Type of job	$T_{arrival}$	Task arrival
T_{ct}	Computational time	$P_{utility}$	Power utility
T_{pd}	Propagation delay	T_{Time}	Tuple time
Fg_{list}	Lists of fog servers	$v_{l, data}$	Volume of data
F_i	Single fog server	d	Length of physical link
F_{broker}	Broker	s	Propagation speed
T_{pt}	Total number of tuples	Σ	End to end delay
F_c	Fog Cache	Q_i	Queuing time
$F_{i, c}$	Computational time of selected fog server	T_d	Transmission delay

4 | PROPOSED METHODOLOGY

In this section, the system model for CBA-based resource allocation for fog computing is shown as shown in Figure 2. In fog computing, job scheduling problem focus at assigning sets of jobs to fog nodes located at the edge of the network in a way to minimize CPU execution time and latency. Table 1 shows the list of symbols and notations that are used throughout the paper.

Let T_{list} denote a set of jobs requested by the IoT devices. The job $T_i \in T_{list}$ is a job requested by the single user and have parameters such as job type (T_{type}), computational time (T_{ct}), and propagation delay (T_{pd}). The type of the job can be small textual, bulk data, location-based, large multimedia, and medical data. Depending on the T_{type} the T_{ct} of the job is computed, which is the processing time of the job. The T_{ct} may vary based on the T_{type} of the job. The T_{pd} is the transmission delay from source to destination. As the processing capabilities of the fog are limited, the T_{type} identify if the job is executed on the fog node or on the cloud, represented as C_{type} and F_{ls} , respectively. If the job is large, then the job is executed on the cloud, otherwise it is executed on an optimal fog node.

The services that are requested by the IoT devices, represented as S_{type} can be categorized as computational, storage, or a combination of both. Fog broker (F_{broker}) serves an an entity that acts as a service provider between the IoT and the fog.²³ We modeled (F_{broker}) as a global gateway that will optimally dispatch the requests between the fog and the cloud. It stores all information about fog servers and also has the information about the current job that is going to be executed.

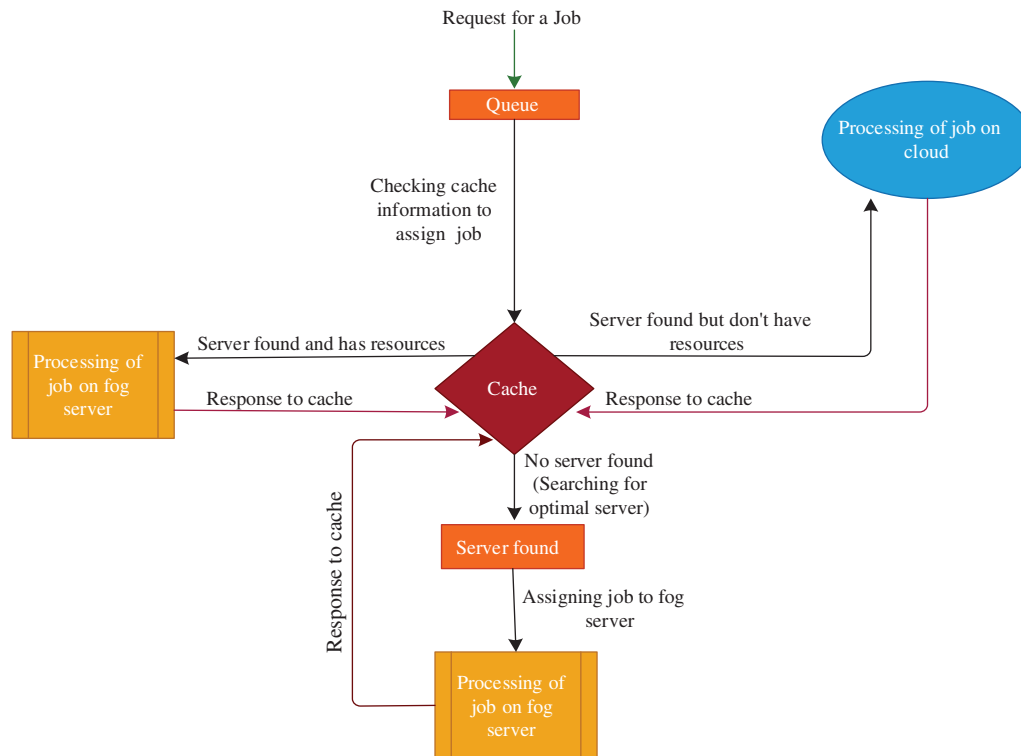


FIGURE 3 Flow diagram of cache-based approach mechanism

The F_{glist} represents the list of fog nodes and F_i is an individual fog node, where $F_i \in F_{glist}$ and $i = \{1, \dots, n\}$ (n is the total number of nodes). We have adopted first in first out (FIFO) methodology to move the jobs from the queue to the cache (F_c). Once the job requests are in the queue, the information from the F_c is used to dispatch them to the respective fog nodes. A queue is maintained to store the jobs, when number of jobs arrive simultaneously. Using FIFO for the jobs in queue, we select nearest fog server and cache the job type, fog server, arrival time, time to leave, and internal processing time.

Once the jobs are dispatched and received at the broker, it performs an initial check to identify if the selected fog server has the required power and resources to execute the job. If the fog nodes has the resources, then its is executed on it. Otherwise, a new node is selected. If no node is available, then the job is moved to the cloud.

The proposed algorithm implements the cache on the smart gateway. Smart gateway helps in connecting IoT devices to the network and efficient utilization of the cloud by setting up the time and type of data to be transferred over the network.^{24,25} The objective of the algorithm is to perform different tasks like preprocessing, filtration of data and reconstruct it in more useful ways, moving data to cloud to provide QoS, keeping a check on delay, execution time, energy consumption. Overall the cache is used in algorithm to expediate the scheduling process and to optimally select the fog nodes.

4.1 | Cache-based job scheduling mechanism

This section presents the proposed CBA scheduling that mainly focuses on designing a centralized mechanism for content delivery by introducing cache-based processing. The devised algorithm with all the functionalities is shown in Algorithm 1. Initially, when a user request for a T_i , the request will go to cache inside the smart gateway followed by the queue (line 1-6 in Algorithm 1) as shown in the flow diagram in Figure 3. After receiving a job request, the proposed algorithm will execute in two steps as follow:

4.1.1 | Execution of job on fog nodes

When job request T_i arrived at the queue, the algorithm first search for the Fog server F_i from the list of available fog servers F_{list} that matches with the requirements of the T_i . If it finds the optimal server for the job, it will assign the job

to that F_i from F_{list} (line 7-14 in Algorithm 1). Next, the type of job and address of the F_i to which the job is assigned is stored in the cache, to save time when a similar type of T_i is requested again (line 17-41 in Algorithm 1). Once done, the algorithm takes the next T_i from T_{list} , which is first checked within F_c whether this type of T_i is executed or not. If yes, then it will send the request to the specified F_i . If the F_i is already fully utilized, then it will make a new entry in F_c and repeat the same procedure (line 42-50 in Algorithm 1). After the execution of job on fog servers the CBA will respond the type of the job and server address to cache to save time in future execution.

4.1.2 | Execution of jobs on cloud

If a new job request arrive and all the F_i are fully utilized, then this request is sent to the cloud data center D_c for processing. Similarly, it also depends on the size of the job, the data or compute intensive jobs are sent to the cloud instead of serving at fog servers as D_c have more processing and storage resources available (line 58-68 in Algorithm 1). After the processing of job CBA will respond backed to cache in order to update the cache information for the future use. The flow diagram of proposed approach is given in Figure 3.

Algorithm 1. CBA scheduling algorithm

Inputs are:

$T_{\text{list}}, F_{\text{glist}}, F_c, Ct = \{\text{true}, \text{false}\}, St = \{S1, S2, S3\}$.

1. **for each:** ($T_i \in T_{\text{list}}$) **do**
2. create T_{Times} class object
3. Set job arrival time
4. Add T_{Time} to T_{list}
5. Checking server
6. compatibility for a job
7. **if** ($F_i \text{Haspower}() \ \&\& \ \text{Hasresources}()$)
8. {
9. **then**
10. create T_{Time} class object
11. set job arrival time
12. creating list of job arrival time
13. Item is not served to cloud
14. **} end if**
15. **end for**
16. **end procedure**
17. PROCEDURE FOGCACHE ($T_{\text{PT}}, F_{\text{GLIST}}$)
18. **if** $F_c.D_{\text{Type}} == F_{\text{broker}.D_{\text{Type}}}.count > 0$
19. {
20. $F_{\text{broker}.fg_i} = \text{haspower}() \ \&\& \ \text{hasresources}$
21. $F_c.Where(x => x.D_{\text{Type}} ==$
22. $F_{\text{Broker}.T_i.D_{\text{Type}}}.OrderBy(x => x.I_{\text{pt}}).OrderBy(x => x.L_i.P_t).get();$
23. }
24. **Else**
25. {
26. $F_{\text{broker}.F_i} = F_{\text{broker}.haspower() \ \&\& \ \text{hasresources}$
27. }
28. **if** ($F_{\text{broker}.F_i}$ is Not null)
29. {CREATE F_{TIMES} AS FOGTIME
30. set job arrival time
31. $F_{\text{Name}} = F_{g_i}$
32. $T_{\text{Name}} = T_i$

```

33. List of jobs having same sources &
34. destination
35.  $T_i$ .IsServed = true;
36.  $T_i$ .IsServerFound = true;
37. Setting endtime in milliseconds
38. create  $T_{Time}$  class object =>  $T_{Time}$ 
39. set departure time
40. } NAME =  $T_I$ 
41. end if
42. if (FogSimulator.IsCreateCache) then
43. {Add elements in list of  $F_c$ 
44. Set  $D_{Type} = T_i.D_{Type}$ 
45. FogServer =  $F_{broker}.Fg_i.ID$ 
46.  $I_{pt} = T_i.I_{pt}$ 
47. TupleGuid =  $T_i.ID$ 
48.  $l_i = L_i$ 
49. }
50. end if
51. %Setting Fog consumption time
52. Set  $F_{Time}.Consumption = P_{Utility}.Consumption(F_{Broker}.F_i, (ttime_{ms} - (InitTime_{ms})), ttime_{ms}, T_i)$ 
53.  $F_{Time}.FreeTime = T_{Time}.T_i$  Departure
54.  $F_{broker}.Fg_i.ReleasePower(F_{broker}.F_i, T_i)$ 
55. Else
56. {Log (“missed by fog”);
57. }
58. if  $C_i == 1$ ) then
59. if (ServedByCloud [tuple, false,  $S_i, D_c$ ]) then
60. {
61. Set  $T_i$ .IsReversed = true
62.  $T_i$ .IsCloudServed = true
63.  $T_i$ .IsServedByFC_Cloud = true
64. }
65. end if
66. else
67. { Set  $T_i$ .IsCloudServed = false
68.  $T_i$ .IsReversed = true
    }
    end if
    Tuple =  $T_i$ 
end FOGCACHE FUNCTION

```

5 | PERFORMANCE EVALUATION

The proposed scheme is evaluated in a C#-based fog computing simulation environment. The simulation setup provides the necessary networking infrastructure, IoT and Fog nodes, and cloud data centers. The simulation platform is motivated by CloudSim²⁶ and iFogSim²⁷ simulators. It provides all the primary and advanced features to simulate IoT-based fog computing environment. We have considered the different types of jobs generated from dumb objects,²⁸ nodes, sensors, mobile, and actuators²⁹ of having datatype small textual, bulk, location-based, large multimedia and medical data. In the experiments, we have considered 30 000 jobs (more details about the jobs can be obtained from Reference 30). There are nine heterogeneous fog servers that serve all the jobs. Fog servers and IoT devices are geographically distributed. IoT devices and fog servers are randomly deployed. The distance between fog servers is in kilometers. Each IoT device is associated with its nearest fog server which is further linked with the cloud datacenters that contain the machines with

high computational capacity and power. The proposed policy is compared with the FCFS and SJF scheduling policies. We used the following metrics to evaluate the performance of our proposed approach.

- propagation delay/latency;
- execution time;
- processing delay; and
- power consumption

Based on the above metrics, the results of the CBA are compared with FCFS and SJF algorithms. In this section, we evaluated our Cache-based technique against FCFS and SJF algorithm to show the effectiveness of CBA for job scheduling in fog computing. The execution of T_i on F_i depends on the T_{type} , T_{ct} , and T_{pd} . So, the execution time, denoted by T_{exec} , is calculated for fog servers as:

$$T_{\text{exec}} = T_{\text{type}} + T_{\text{ct}} + T_{\text{pd}}. \quad (1)$$

Equation (1) mentioned here is commonly used in the literature such as References 31,32. The computational time of job T_i in fog server, represented as $F_{i,c}$ is calculated as:

$$F_{i,c} = \frac{v_{t,\text{data}}}{\sum_{p=1}^{P_{\text{max}}} r_{t,s}^j}, \quad (2)$$

where $r_{t,s}^j$ represents resources that are allocated by fog node for T_i during the period p , where P_{max} is the maximum time that a job can maintain, and $v_{t,\text{data}}$ represents the volume of data that need to be processed. Furthermore, the end-to-end delay is given as:

$$\Sigma = Q_t + T_d + T_{\text{pd}}, \quad (3)$$

where Σ represents end-to-end delay, which is computed as summation of Queuing time (Q_t), that is, time taken by job to wait in the queue until it can be executed, transmission delay (T_d), and propagational delay (T_{pd}). The T_{pd} can be calculated in the similar manner as³³:

$$T_{\text{pd}} = \frac{d}{s}, \quad (4)$$

where d is the length of physical link and s is the propagation speed in the medium. Transmission delay can be defined as how long it takes to get all the packets into the wire in the first place. Transmission delay is calculated as follow:

$$T_d = \frac{L}{R}, \quad (5)$$

where L is the length of the job/tuple and R is the link bandwidth in bits per second. *Execution time* means the time required to complete simulation, which means the scheduling of all the jobs. Three different algorithms are run on the same dataset and on the same machine to evaluate the results. In Figure 4, it is shown that our proposed model takes less time in simulation in comparison to two other state-of-the-art algorithms. Cache takes 4482.928 seconds to complete the simulation as compared with FCFS and SJF which took 6843.772 seconds and 7280 seconds, respectively. In Figure 4, the reason for less execution time is that cache information is being used to assign the job to fog server. Jobs do not have to wait for long in the queue if all servers are busy then it starts sending jobs to the cloud for execution.

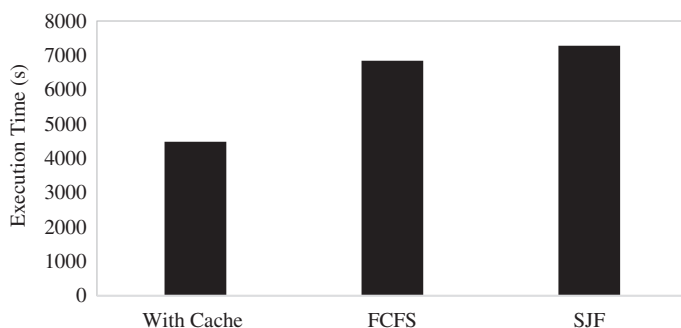


FIGURE 4 Execution time

FIGURE 5 Propagation delay

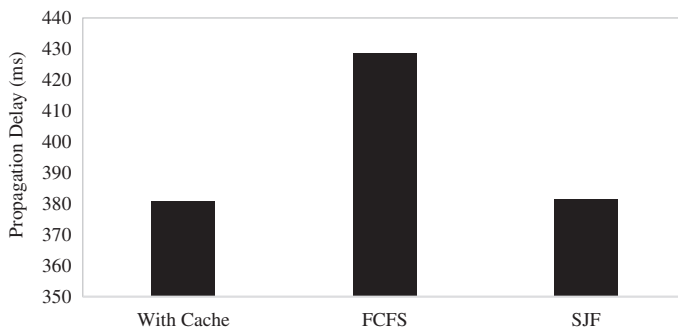


FIGURE 6 Total average internal processing delay

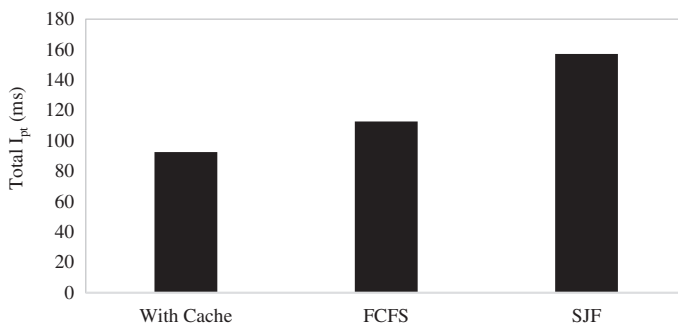
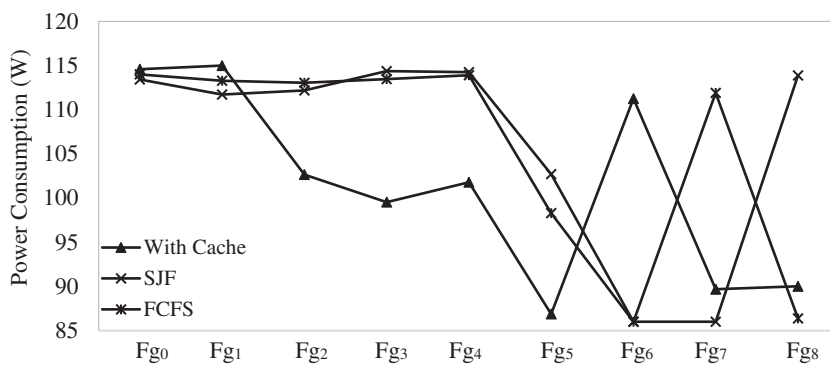


FIGURE 7 Server power consumption



Propagation delay is the time taken by a job from the source to reach the destination. We calculated the average propagation delay of cache along with FCFS and SJF algorithm as shown in Figure 5. Time taken by job request is calculated to reach from source to the destination server. Average propagation delay taken by our proposed model is less than the other two algorithms as shown in Figure 5. The average propagation delay took by each job to reach to fog server while using cache approach is 380.853 ms. FCFS takes 428.594 and SJF takes 381.503 ms, respectively. In fact, the affective distribution of jobs to the available resources results in better bandwidth utilization that ultimately reduces the congestion on the network and results in lower propagation delay. Similarly, the purpose of our proposed algorithm is to utilize the fog resources as maximum as possible that also has an impact on the propagation delay—the jobs sent to cloud results in higher delays.

In our evaluation test, the delay is measured in milliseconds, and it is the time taken by a processing element to execute the job. Figure 6 presents the average internal processing delay of the servers. The cache takes less time when compared with FCFS and SJF. Cache takes 92.567 ms whereas FCFS takes 112.648 and SJF takes 157.106 ms, respectively.

The amount of processing performed by the servers is directly proportional to the power consumed by the servers. In Figure 7, we have depicted the power consumption of fog servers in Watts (W). The x-axis presents the servers ID and y-axis exhibits their corresponding power consumption. Figure 7 shows the trend of power consumed in the fog resources. The power consumption has higher values throughout the time because of its two types—static and dynamic power consumption. Static power consumption refers to the power required for the working of electronic peripherals of fog servers when it is turned on and there is no load on it. Afterward, the power consumption is proportional to its

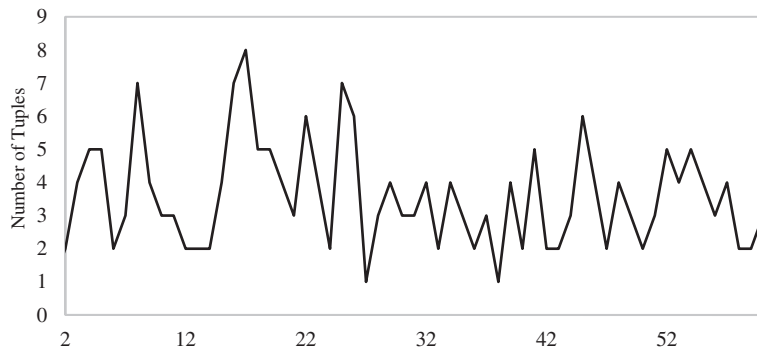


FIGURE 8 Tuple execution

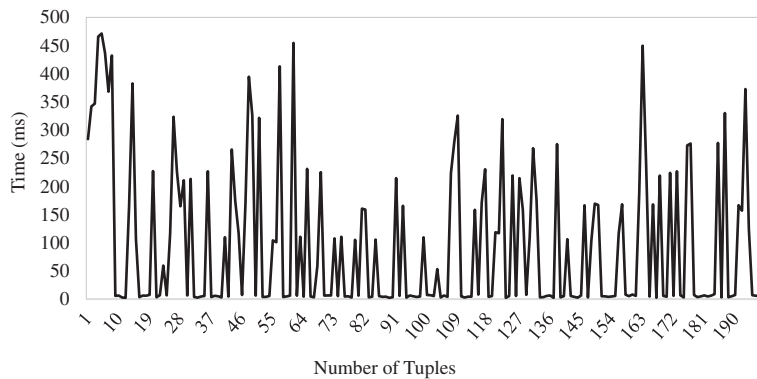


FIGURE 9 End-to-end delay of tuples

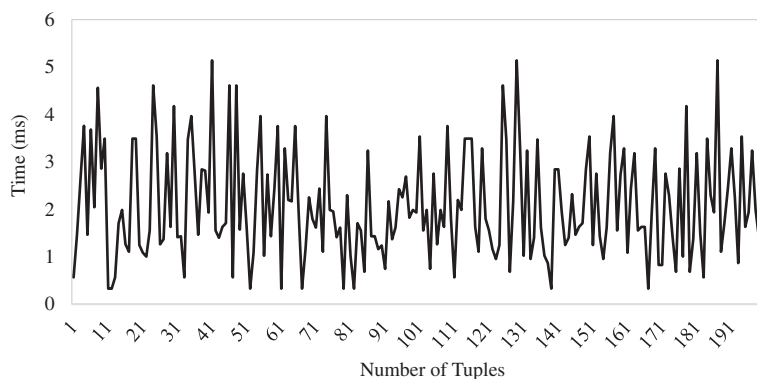


FIGURE 10 Network propagation time

utilization—called dynamic power consumption. The utilization of the fog servers depends on the millions of instructions per second required by the jobs at fog server for its computation and also the frequency of jobs has an impact on the utilization of fog resources. Here, we refer the fog servers to fog resources. The utilization of the fog servers depends on the scheduling policy. Figure 7 shows that fog servers are utilized arbitrarily. It can be seen that some servers have more consumption that represents that they capture more traffic than the others. When fog resources are saturated, the incoming requests are forwarded to the cloud. Similarly, it also depends on the job size, the big jobs are sent to the cloud instead of serving at fog servers.

Tuple execution time is the average time taken by a tuple to execute. In Figure 8 it is shown that how tuples are executed in each second. As in Figure 8 in the first second, only one job is executed, and in 17 seconds, eight jobs are executed. Each tuple has an associated CPU and network cost to process it.³¹

Figure 9 shows the average end-to-end delay experienced by tuple to execute on fog and cloud using the CBA mechanism. It can be noticed that the average end-to-end tuple delay falls below as data is processed near to source and cache information is used due to which tuple does not have to wait for long in the queue. The reasons for peak values in the graph is due to the execution of jobs on the cloud and, the distance between the job source and the fog server. In Figure 10 the time taken by each tuple to travel from sender to destination is depicted.

Based on the above results the overall performance increased. One disadvantage of this method is that initially, server utilization is maximum which can be a future work to minimize it. The overall performance increased, latency and propagation delay decreased, and QoS achieved.

6 | CONCLUSION AND FUTURE WORK

Since the inception of IoT, edge and fog computing paradigm have undergone an enormous evolution in a way that they can be used. More and more new applications, such as face recognition, augmented reality, online interactive gaming, and natural language processing are emerging and attracting the researcher to explore methods to enable computing near device level. However, such applications are generally data intensive or compute intensive, which demands high resource and energy consumption. Therefore, enable sophisticated computing algorithms at fog or edge node which are resource constraint devices is a challenging task. In this paper, the job scheduling problem in the fog computing environment for the efficient execution of tasks requested by end user devices is explored. The proposed algorithm integrates cache in the smart gateway and proposes a scheduling scheme that decreases the execution time, propagation delay, and internal processing time of the jobs that are being requested. To handle job scheduling the proposed algorithm uses FCFS policy for a queue. The performance of the proposed algorithm is compared with the traditional FCFS and SJF policies, and the results showed that our approach is more optimized and yields reduction of execution time, latency, processing delays and power consumption by 38%, 11.1%, 6%, and 17.8%, respectively, as compared to the FCFS and SJF policies. In the future, the aim is to extend the proposed algorithm to minimize the server execution as it is very high at the start of the algorithm and to explore SJF policy in a CBA-based resource allocation scheme for fog computing.

ORCID

Saif U. R. Malik  <https://orcid.org/0000-0001-8195-1630>

Haris Pervaiz  <https://orcid.org/0000-0002-8364-4682>

REFERENCES

1. Malik H, Kandler N, Alam MM, Annus I, Le Moullec Y, Kuusik A. Evaluation of low power wide area network technologies for smart urban drainage systems. Paper presented at: 2018 IEEE International Conference on Environmental Engineering (EE); 2018; Milan, Italy:1–5.
2. Alam MM, Malik H, Khan MI, Pardy T, Kuusik A, Le Moullec Y. A survey on the roles of communication technologies in IoT-based personalized healthcare applications. *IEEE Access*. 2018;6:36611–36631.
3. Saleem A et al. FESDA: fog-enabled secure data aggregation in smart grid IoT network. *IEEE Internet of Things J*. 2019.
4. Datta SK, Bonnet C, Haerri J. Fog computing architecture to enable consumer centric Internet of Things services. Paper presented at: 2015 International Symposium on Consumer Electronics (ISCE); 2015; Madrid:1–2.
5. Luan TH, Gao L, Li Z, Xiang Y, Wei G, Sun L. Fog computing: focusing on mobile users at the edge. *Comput Sci*. 2015;1–11.
6. Yi S, Li C, Li Q. A survey of fog computing: concepts, applications and issues. Paper presented at: Proceedings of 2015 Workshop on Mobile Big Data (MBD'15); June 2015:37–43.
7. Malik SUR, Khan SU, Ewen SJ, et al. Performance analysis of data intensive cloud systems based on data management and replication: a survey. *Distrib Parallel Databases*. 2016;34(2):179–215.
8. Idachaba U, Wang F. A community-based cloud computing caching service. Paper presented at: 2015 IEEE International Congress on Big Data; 2015; NewYork, NY:559–566.
9. Ruia A, Casey CJ, Saha S, Sprintson A. Flowcache: A cache-based approach for improving SDN scalability. Paper presented at: 2016 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPS); 2016; San Francisco, CA:610–615.
10. Assila B, Kobbane A, El Koutbi M. A many-to-one matching game approach to achieve low-latency exploiting fogs and caching. Paper presented at: 2018 9th IFIP International Conference on New Technologies, Mobility and Security (NTMS); 2018; Paris, France:1–2.
11. Aazam M, Huh E-N. Fog computing and smart gateway based communication for cloud of things. Paper presented at: in Proceedings of International Conference on Future Internet Things Cloud (FiCloud); 2014:464–470.
12. Krishnapriya S, Joby PP. QoS aware resource scheduling in internet of things-cloud environment. *Int J Sci Eng Res*. 2015;6(4):294–297.
13. Bitam S, Zeadally S, Mellouk A. Fog computing job scheduling optimization based on bees swarm. *Enterprise Inform Syst*. 2017;12(4):1–25.
14. Verma M, Bhardwaj N, Yadav A. Real time efficient scheduling algorithm for load balancing in fog computing environment. *Int J Inform Technol Comput Sci*. 2016;4:1–10. <https://doi.org/10.5815/ijitcs.2016.04.01>.
15. Intharawijitr K, Iida K, Koga H. Analysis of fog model considering computing and communication latency in 5G cellular networks. Paper presented at: 2016 IEEE International Conference on Pervasive Computing and Communication Workshops (PerCom Workshops); 2016; Sydney, NSW:1–4.
16. Ningning S, Chao G, Xingshuo A, Qiang Z. Fog computing dynamic load balancing mechanism based on graph repartitioning. *China Commun*. 2016;13(3):156–164.

17. Yin L, Luo J, Luo H. Tasks scheduling and resource allocation in fog computing based on containers for smart manufacturing. *IEEE Trans Ind Inform*. 2018;14(10):4712-4721.
18. Choudhari T, Moh M, Moh TS. Prioritized task scheduling in fog computing. Paper presented at: Proceedings of the ACMSE 2018 Conference; March 29–31, 2018; Richmond, Kentucky:1–8
19. Idachaba U, Wang F. A community-based cloud computing caching service. 2015 IEEE International Congress on Big Data; 2015; NewYork, NY:559–566.
20. Jia A, Han G, Wang H, Wang F. Cost-aware cache replacement policy in shared last-level cache for hybrid memory based fog computing. *Enterprise Inform Syst*. 2018;12(4):435-451.
21. Aazam M, Zeadally S, Harras KA. Fog computing architecture, evaluation, and future research directions. *IEEE Commun Mag*. 2018;56(5):46-52.
22. Hu P, Dhelim S, Ning H, Qi T. Survey on fog computing: architecture, key technologies, applications and open issues. *J Netw Comput Appl*. 2017;98:27-42.
23. Gedeon J, Meurisch C, Bhat D, Stein M, Wang L, Mühlhäuser M. Router-based brokering for surrogate discovery in edge computing. 2017 IEEE 37th International Conference on Distributed Computing Systems Workshops (ICDCSW); 2017; Atlanta, Georgia:145–150.
24. Dastjerdi AV, Gupta H, Calheiros RN, Ghosh SK, Buyya R. Fog computing: principles, architectures, and applications. *Internet of Things*. Burlington, MA: Morgan Kaufmann; 2016:61-75.
25. Aazam M, Hung PP, Huh E. Smart gateway based communication for cloud of things. Paper presented at: 2014 IEEE Ninth International Conference on Intelligent Sensors, Sensor Networks and Information Processing (ISSNIP); 2014; Singapore:1-6.
26. Calheiros RN, Ranjan R, Beloglazov A, Buyya R. CloudSim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms. *Softw Pract Exp*. 2011;41:23-50.
27. Gupta H, Vahid Dastjerdi A, Ghosh SK, Buyya R. iFogSim: a toolkit for modeling and simulation of resource management techniques in the internet of things, edge and fog computing environments. *Softw Pract Exp*. 2017;47(9):1275-1296.
28. Aazam M, Huh E. Fog computing micro datacenter based dynamic resource estimation and pricing model for IoT. Paper presented at: 2015 IEEE 29th International Conference on Advanced Information Networking and Applications; 2015; Gwangju:687–694.
29. Lee W, Nam K, Roh H, Kim S. A gateway-based fog computing architecture for wireless sensors and actuator networks. Paper presented at: 2016 18th International Conference on Advanced Communication Technology (ICACT); 2016; Pyeongchang:1–1.
30. GITHUB. Iot-compute-dataset. 2019. <https://github.com/saifulislamPhD/IoT-647ComputeDataset>. Accessed November 22, 2019.
31. Xiao Y, Krunz M. Qoe and power efficiency tradeoff for fog computing networks with fog node cooperation. Paper presented at: Proceedings of IEEE Conference on Computer Communications; 2017:1–9.
32. Pang A, Chung W, Chiu T, Zhang J. Latency-driven cooperative task computing in multi-user fog-radio access networks. Paper presented at: 2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS); 2017; Atlanta, GA:615–624.
33. Akbar S, Malik SUR, Khan SU, Choo R, Anjum A, Ahmad N. A game-based thermal-aware resource allocation strategy for data centers. *IEEE Trans Cloud Comput*. 2019. <https://doi.org/10.1109/TCC.2019.2899310>.

How to cite this article: Khan OA, Malik SUR, Baig FM, et al. A cache-based approach toward improved scheduling in fog computing. *Softw: Pract Exper*. 2021;51:2360–2372. <https://doi.org/10.1002/spe.2824>