

# A Novel Data Placement and Retrieval Service for Cooperative Edge Clouds

Junjie Xie, Chen Qian, Deke Guo, Xin Li, Ge Wang, Honghui Chen

**Abstract**—Mobile edge computing is a new paradigm in which the computing and storage resources are placed at the edge of the Internet. Data placement and retrieval are fundamental services of mobile edge computing when a network of edge clouds collaboratively provide data services. These services require short-latency and low-overhead implementation in network and computing devices and load balance on edge clouds. However existing methods such as distributed hash tables (DHTs) are not enough to achieve efficient data placement and retrieval services for cooperative edge clouds. This paper presents GRED, a novel data placement and retrieval service for mobile edge computing, which is efficient in not only the load balance but also routing path lengths and forwarding table sizes. GRED utilizes the programmable switches to support a virtual-space based DHT with only one overlay hop. Data location can be easily implemented on top of the GRED by associating a virtual position with each data by hashing, and storing the data at the edge server connected to the switch whose position is the nearest to the position of the data in the virtual space. We implement GRED in a P4 prototype, which provides a simple and efficient solution. Results from theoretical analysis, simulations, and experiments show that GRED can efficiently balance the load of edge clouds, and can fast answer data queries due to its low routing stretch.

**Index Terms**—Data placement, Data retrieval, Cooperative edge clouds, Mobile edge computing.

## 1 INTRODUCTION

CLOUD computing is a common solution to provide resources for computation, storage, and bandwidth to massive mobile computing devices such as those of the Internet of Things (IoT). However, many modern applications such as augmented reality (AR), wearable cognitive assistance, and real time monitoring/control are latency-sensitive and may suffer the long round-trip delay to the Cloud. A recent trend is to offload computing and storage to the network edges so as to enable computation-intensive and latency-critical applications. This technology, called *Edge Computing* [1], [2], has been proposed to shift computing and storage capacities from the remote Cloud to the network edge in close proximity to mobile devices, sensors, and end users. Edge computing promises the dramatic reduction in network latency and traffic volume, tackling the key challenges for materializing the 5G vision. The edge of the Internet offers ideal placement for low-latency offload infrastructure to support emerging applications. Terms such as ‘cloudlets’, ‘micro data centers’, and ‘fog computing’ have been used in the literature to refer to similar edge-located services [3] [4].

In edge computing, each edge cloud consisting of mul-

iple edge servers has certain computation and storage resources, and this provides a chance to offload part of the workload from the remote Cloud. On the one hand, the edge clouds would cache the data from the remote Cloud [5]. On the other hand, the edge users would store their application data in edge clouds. Meanwhile, sharing the computation and storage resources among those edge clouds can balance the uneven distribution of the computation and storage workloads and capabilities over edge clouds. Unlike Cloud datacenters, edge clouds are usually geographically distributed and have heterogeneous computation and storage capacities [1]. Those ad hoc-like connected edge clouds provide the opportunity for stakeholders to share and cooperate data and resources where the edge clouds of multiple stakeholders are geographically distributed [6]. Always offloading the data and computation of users at the closest edge cloud may not be a valid solution because 1) the user may be mobile and 2) one edge cloud has limited resources. Hence we consider a large number of edge clouds in an interconnected edge network that collaboratively serve the resource pool of storage and computation offloading for users.

A core operation for the cooperative edge clouds is to support the efficient *data placement and retrieval* when multiple edge clouds work together [7]. In this work, we define “data placement” as the process of delivering a given data item to an edge server for storage and “data retrieval” as the process of finding the storage server of a given data item and requesting the server to deliver the data to a user. Hence they are essentially *overlay services with network-layer implementation*. More importantly, the data placement and retrieval services will efficiently support a large number of upper-layer applications in the environment of mobile edge computing. The data problems in edge computing are very fundamental and urgent. However, the data placement and

- Junjie Xie is with the Institute of Systems Engineering, AMS, PLA, Beijing, 100141, P.R. China. E-mail: xiejunjie06@gmail.com.
- Deke Guo and Honghui Chen are with the Science and Technology Laboratory on Information Systems Engineering, National University of Defense Technology, Changsha Hunan, 410073, China. E-mail: {guodeke, chh0808}@gmail.com.
- Chen Qian and Xin Li are with the Department of Computer Science and Engineering, University of California Santa Cruz, CA 95064, USA. E-mail: cqian12@ucsc.edu.
- Ge Wang is with the Department of Computer Science and Engineering, Xi'an Jiaotong University, China. Email: wangge@stu.xjtu.edu.cn.

retrieval services in mobile edge computing face at least two challenges. First, these services should have short-latency and low-overhead implementation on the user side and network routers/switches. For example, it is impractical to maintain a complete index of all data-to-location mappings at an edge device or inside a router. Second, achieving load balance among edge clouds is very important, which requires that no server should be overloaded when there is the available resource on other servers. The limited and heterogeneous computation and storage capacities of different edge clouds further complicate the problem.

To solve these problems, we propose short-latency and low-overhead data placement and retrieval services for cooperative edge clouds, called Greedy Routing for Edge Data (GRED). GRED includes two innovative ideas. First, GRED supports a DHT of edge clouds with only one overlay hop. Second, GRED utilizes the Software Defined Networking (SDN) paradigm [8], [9], [10] to implement efficient routing support of the one-hop DHT on programmable switches<sup>1</sup>. In particular, the SDN controller maintains a virtual space. Switches and data items are mapped to different positions in the space according to their IDs. The data will be stored in an edge server connected to the switch whose position is nearest to the data position in the virtual space.

**GRED is efficient in terms of both routing path lengths and forwarding table sizes.** Each data placement/retrieval request in GRED only needs one overlay hop. In detail, in the control plane of GRED, we design a virtual space construction algorithm to assign the switches to the points in the virtual space, such that the Euclidean distance between two switches is proportional to their network distance. It is proved that under this circumstance, the routing stretch of the network can be optimized. Furthermore, to achieve the load balance among edge clouds, we further optimize the switches' positions considering that the data is stored in the network based on their positions in the virtual space.

Meanwhile, to minimize the forwarding table size, the data plane of GRED does not need a new flow entry for every placement/retrieval request. Instead, the data plane performs greedy forwarding based on the next-hop switch's position determined by greedy forwarding, which is implemented in P4 [11], a programmable data plane development tool. Hence the forwarding table size is independent of the network size and the number of flows in the network. However, some edge servers with low capacities would be overloaded when switches connects to the heterogeneous edge servers. To solve this problem, we further design the extended-GRED protocol to enlarge the management range of switches. More precisely, the control plane would update the flow entries of the related switches, which can redirect the data to an edge server that still has the remaining capacity. Although the extended-GRED protocol would incur a little higher routing stretch, it can completely eliminate the overload of edge servers. We conducted extensive experiments, using both P4 implementation and simulations, to evaluate the performance of GRED. Theoretical analysis shows the correctness and efficiency of GRED. Experimental results show that GRED uses <30% routing cost and

achieves better load balance among edge clouds compared to using Chord [12], a well-known DHT.

We summarize our contributions as follows.

- 1) We study the data placement and retrieval problem among edge-clouds in the edge computing environment with the aim to provide low-latency and low-overhead data services.
- 2) We propose a greedy data routing approach to efficiently provide a feasible routing path with shorter lengths and smaller table sizes, comparing with other DHTs-based approaches, i.e. Chord.
- 3) We evaluate the performance of the GRED protocol using an implementation on our testbed and massive simulations. The results of those experiments show the efficiency and effectiveness of the GRED protocol.

The rest of this paper is organized as follows. In section 2, we introduce the motivation and preliminaries of this paper. Section 3 presents the system overview. In Section 4, we describe the virtual position construction in the control plane, which is the base of the GRED. Section 5 details the placement and retrieval mechanism. We discuss the network dynamic and the data copies in Section 6. In Section 7 we evaluate the performance of the GRED. We introduce the related work and conclude this paper in Section 8 and Section 9, respectively.

## 2 MOTIVATION AND PRELIMINARIES

In this section, we first present the motivation of this paper and then introduce some preliminaries about the Delaunay Triangulation (DT) graph [13].

### 2.1 Motivation

A recent trend is to offload computing and storage to the network edges so as to enable computation-intensive and latency-critical applications. Edge computing promises dramatic reduction in latency and energy consumption, tackling the key challenges for materializing 5G vision. The promised gains of edge computing have motivated extensive efforts in both academia and industry on developing the technology [14]. This work focuses on the core function in edge computing: data placement and retrieval, which provide efficient support for a large number of emerging applications. First, these services should have the efficient implementation on the user side and network routers/switches. The efficiency discussed in this work aims at both network efficiency such as short routing path lengths that imply short latency and low bandwidth cost, and forwarding efficiency such as small forwarding table size that achieves low infrastructure cost. The second challenge is the load balance when a large amount of data [15] are stored in those edge servers. Load balance requires that no edge server should be overloaded when there is the available resource on other servers. The limited and heterogeneous computation and storage capacities of different edge clouds further complicate the problem.

However, the data placement and retrieval services are confront with two key challenges. First, an unavoidable issue involves the edge computing is the user mobility.

1. Hereafter we use "switches" to denote network forwarding devices for the compatibility to the SDN context, although they can be routers.

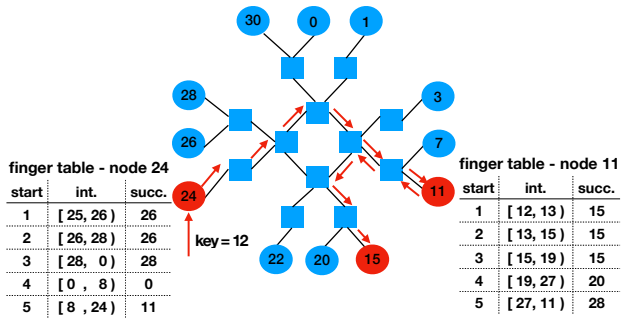


Fig. 1. Finger tables and key locations in DHT-based storage system.

That is, a user hopes to continue to use the data when the user moves from one location to another location. Another scenario is that a user could hope to share the data with other users who are in different locations. Meanwhile, the edge servers need to cache the data from the remote Cloud for the edge users. However, it is very hard to estimate the distribution of potential users. Therefore, the edge clouds need to provide efficient data lookup services for the edge users that access the network from any edge server. To enable those scenarios, the critical problem is to locate the data wherever a user accesses the network.

Another important issue involves the edge computing is the load balance when a large amount of data are stored in those edge nodes. Although the edge servers can conduct the computing and storage tasks, the edge servers are with limited computing and storage resources. Furthermore, the computing and storage tasks could be uneven distributed in those edge clouds. Sharing and cooperating data among multiple edge clouds is an efficient solution to the above problems. However, how to achieve the load balance among multiple edge clouds is another challenge that we need to solve in this paper.

To enable the data placement and retrieval services, we recall that there has been some related work in peer-to-peer (P2P) networks. However, existing approaches in P2P can not meet the low-latency routing requirement in edge computing. In those systems, a data object is associated with a key and each node in the system is responsible for storing a certain range of keys. For example, Chord [12] is a widely used DHT solution for the data storage and lookup in P2P networks. As shown in Fig. 1, an edge network consists of 12 edge servers where each edge server has a unique identifier. The data with the key 12 is stored in server 15 based on the storage principle in Chord [12]. When a user accessing server 24 needs to retrieve the data with the key 12 located in the interval [8, 24), the lookup request is first sent to server 11 based on the finger table of server 24. Note that the finger table indicates the successor node to find a data. Then, server 11 will continue to forward the lookup request to sever 15 based on its finger table. In this case, the path length of the lookup request is 11, which is significantly longer than the shortest path between server 24 and server 15 with only 5 physical hops.

In the DHT-based storage systems, the overlay routing takes up to  $O(\log n)$  overlay hops for  $n$  nodes and each overlay hop may include multiple network-layer hops [16][17]. The main reason is the mismatch between the overlay net-

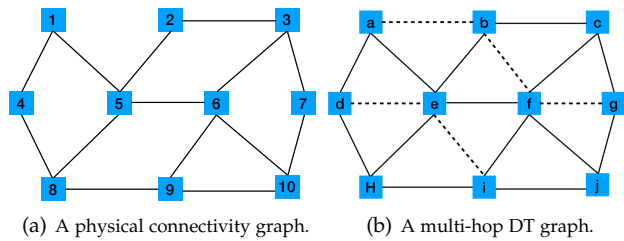


Fig. 2. An illustration of a physical network and the multi-hop DT

work and the physical network. That is, the path length for locating a data item is heavily longer than the shortest path. Furthermore, such mismatch causes a high routing stretch, which results in the long response delay. Note there has been some research proposed to improve the overlay-network mismatch problem [18] [19]. However experimental results show that they are not able to maintain a routing stretch lower than 2 for large networks [13]. Although some work can achieve  $O(1)$  DHT, such as Structured Superpeers [20] and Beehive [21], they need to store a large amount of indexing information or add many data duplicates in edge servers. In addition, the load balance in Chord [12] is not good enough. Although Chord can achieve a better load balance by adding more virtual nodes to each real node, it also increases the routing table space usage and makes the system more complicated. Therefore, in this paper, we look for a better design with low routing stretch and better load balance to enable the data placement and retrieval service for cooperative edge clouds.

## 2.2 Guaranteed delivery on a DT Graph

In our design, each switch does a greedy forwarding. To achieve the guaranteed delivery, a virtual DT graph is maintained in the control plane of the network. Note that greedy routing on an arbitrary graph is prone to the risk of being trapped at a local optimum, i.e., routing stops at a non-destination node that is closer to the destination than any of its neighbors. However, on a DT, it is guaranteed that greedy routing always succeeds to find the node closest to destination  $p$ . For a given set  $P$  of discrete points (called nodes) in a plane is a triangulation  $DT(P)$  such that no point in  $P$  is inside the circumcircle of any triangle in  $DT(P)$ . If two nodes share a DT edge, they are called DT neighbors. One important property of DT is that greedy routing to a destination location  $p$  on a DT graph always stops at a node that is closest to  $p$  among all nodes [22].

Note the main difficulty of maintaining a DT graph in a network of edge nodes is that two DT neighbors may not be connected by a physical link. Hence they cannot directly forward messages between them. For an arbitrary layer-2 network, the MDT [13] protocol was designed for nodes to construct a distributed multi-hop DT graph. As shown in Fig. 2(b), there is a DT graph of 10 nodes in a 2D Euclidean space, and the physical connectivity of those 10 nodes is shown in Fig. 2(a). In Fig. 2(b), Nodes 5 and 1 are both physical neighbors and DT neighbors. However, DT neighbors, nodes 1 and 2, are not connected directly in Fig. 2(a). Hence in a multi-hop DT graph, node 1 transfers packets to node 2 by the multi-hop path  $\{1, 5, 2\}$  in Fig.

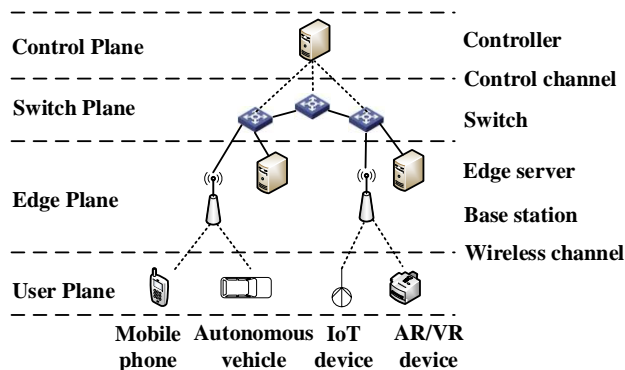


Fig. 3. General architecture of software-defined edge network.

2(a). Therefore, node 2 is called the multi-hop DT neighbor of node 1. For a set of nodes that maintain a correct multi-hop DT, given a destination  $p$ , it is proved that MDT-greedy forwarding always succeeds to find a node that is closest to  $p$ , for nodes located in a Euclidean space (2D, 3D, or a higher dimension) [13].

### 3 SYSTEM OVERVIEW

The GRED protocol specifies how to place a data item and to retrieve it from the edge servers given a data identifier. We design the GRED protocol while utilizing the advantage of software-defined networking [8], which centralizes the network intelligence in the network controller. The switches in the data plane only forward packets according to the related forwarding entries derived from the controller. When we apply the principle of SDN to the edge computing, the network is called a Software-Defined Edge Network (SDEN). As shown in Fig. 3, we define the general hierarchical architecture of an SDEN, which consists of the control plane, the switch plane, the edge plane and the user plane. The user plane includes the mobile users and various edge devices, such as autonomous vehicles and IoT devices. In SDEN, the users access the network by wireless Access Points (APs). Those APs and edge servers are connected to network switches and constitute the edge plane. The switches provide data communication services among edge servers based on the forwarding entries derived from the controller in the control plane [9], [10].

The GRED protocol mainly consists of the functions in the control plane and the switch plane.

**Control plane** associates each switch to a point in the virtual space and computes a DT graph of all points. It then inserts related forwarding entries into switches based on their DT neighboring relationships in the virtual space. It is worth noting that the control plane proactively distributes forwarding entries to switches, which perform greedy forwarding based on the destination position rather than per-flow information. The mechanism can efficiently reduce the load of the control plane and the size of forwarding tables, because the switches can forward data requests based on the pre-installed rules without the interaction of the control plane.

**Switch plane** consists of switches and transfer links. The switch greedily forwards a data request to the correct edge server based on the installed rules. More precisely, the

switch first achieves the data position in the virtual space by hashing the data identifier. Then the switch finds a DT neighbor that is closest to the data position and forwards the packet to it, by either a direct link or a multi-hop path.

When placing a data item to an edge server, the hash value  $H(d)$  of the data identifier  $d$  is firstly calculated. In this paper, we adopt the hash function, *SHA-256* [23], which outputs a 32-byte binary value. Furthermore, the hash value  $H(d)$  is reduced to the scope of the 2D virtual space, which is constructed by the control plane. We only use the last 8 bytes of  $H(d)$  and convert them to two 4-byte binary numbers,  $x$  and  $y$ . We limit so that the coordinate value ranges from 0 to 1 in each dimension. Then, the position of a data in 2D is  $(\frac{x}{2^{32}-1}, \frac{y}{2^{32}-1})$ . The position can be stored in decimal format, using 4 bytes per dimension. Hereafter, for any data identifier,  $d$ , we use  $H(d)$  to represent its position. Last, the data is greedily forwarded to the switch whose position is the nearest to the data position in the virtual space, and further, the switch determines a unique edge server to store the data.

We design the GRED protocol keeping the following three goals in mind:

- 1) *Guaranteed delivery*: Given a data identifier, the GRED forwarding protocol always succeeds to find a switch closest to the data location. Furthermore, the switch determines a unique edge server for the data.
- 2) *Low routing stretch*: The GRED forwarding path is close to the shortest path between the edge server storing the data and the edge server sending the data request wherever the data request is sent.
- 3) *Load balance*: GRED should place data among all cooperative edge servers such that no server is overloaded.

The GRED protocol greedily forwards the data request based on the data position and the switches' positions in the virtual space. Determining the positions of switches is the key to achieve the advantages of GRED. It is because bad virtual positions will result in long routing path and bad load balance among edge servers. Therefore, we first detail the procedure of the virtual position construction in the next section.

## 4 VIRTUAL POSITION CONSTRUCTION

The control plane of GRED first determines the positions of all switches in a virtual 2D Euclidean space, then constructs a multi-hop DT [22] based on those virtual positions. After that, the control plane inserts the forwarding entries into the switches. Then, switches perform greedy forwarding based on those forwarding entries. The key point is to determine the positions of the switches, which affect the routing stretch and the load balance of the GRED protocol.

### 4.1 Calculating the coordinates of switches

Note that the network topology and state can be obtained in the control plane by collecting switch, port, link, and host information [10][24]. Then, the control plane can compute the shortest path matrix between switches. To ensure the low routing stretch of greedy routing, it is required that

the Euclidean distance of two switches in the virtual space is proportional to their network distance, which is called greedy network embedding [25]. Therefore, the key challenge is how to achieve the coordinate matrix of  $n$  points where the shortest path lengths between  $n$  switches can be indirectly reflected by the distances between  $n$  points in the virtual space. To achieve this goal, we design *Scoord* algorithm to calculate the switches' coordinates in the virtual space.

The *Scoord* algorithm utilizes the Multidimensional Scaling (MDS) technique [26]. The MDS aims to place each object in  $m$ -dimensional space such that the between-object distances in the distance matrix are preserved as well as possible in the Euclidean distances in the space. Each object is then assigned coordinates in the  $m$  dimensions. The number of dimensions  $m$  of MDS can exceed 2 and is specified in advance. Choosing  $m=2$  optimizes the object locations for a 2D Euclidean space. Inspired by the MDS, we design the *Scoord* algorithm to calculate the positions of switches in the virtual space while preserving the network distances between switches. The *Scoord* algorithm takes an input matrix giving network distances between pairs of switches, which is known to the control plane.

In the control plane, the shortest path matrix  $\Gamma$  between switches is first calculated. The *Scoord* algorithm utilizes the fact that the coordinate matrix can be derived by eigenvalue decomposition from  $B=QQ'$  where matrix  $B$  can be computed from the distance matrix  $\Gamma$  [27]. Then, the matrix  $Q$  can be uniquely determined by matrix  $B$ . Therefore, the *Scoord* algorithm first constructs the scalar product matrix  $B$  by multiplying the squared distance matrix  $\Gamma^{(2)}$  with the matrix  $\Theta=I-\frac{1}{n}\Delta$ , where  $\Delta$  is the squared matrix with all elements are 1, and  $n$  is the number of switches. Then,  $B=-\frac{1}{2}\Theta\Gamma^{(2)}\Theta$ . This procedure is called double centering [28]. Furthermore, the  $m$  largest eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_m$  and the corresponding eigenvectors  $e_1, e_2, \dots, e_m$  of the matrix  $B$  is determined, where  $m$  is the number of dimensions. Last, the coordinates of the switches  $Q=E_m\Lambda_m^{1/2}$  are achieved, where  $E_m$  is the matrix of  $m$  eigenvectors and  $\Lambda_m$  is the diagonal matrix of  $m$  eigenvalues of the matrix  $B$ , respectively. After that, each switch will be assigned a coordinate in the virtual space from the coordinate matrix  $Q$ . When the control plane maintains a 2D Euclidean space,  $\Lambda_m$  is the diagonal matrix of 2 largest eigenvalues of the matrix  $B$ , and  $E_m$  is the matrix of 2 corresponding eigenvectors.

## 4.2 Refining the positions of switches

One potential problem is that the *Scoord* algorithm determines the positions of switches without considering the load balance among edge nodes. Fig. 4(a) shows a Voronoi Diagram [29] of 10 crosses where each cross is associated with a Voronoi cell. In each cell, the distance from a point to the corresponding cross in the same region is not greater than its distance to the other crosses in the diagram. Recall that a data item is stored at an edge server connected to the switch whose position is nearest to the data position in the virtual space. Assume that the switches are located in those crosses in Fig. 4(a). Then, when data items are mapped into the whole space evenly, it is obvious that there would be load imbalance among those switches because

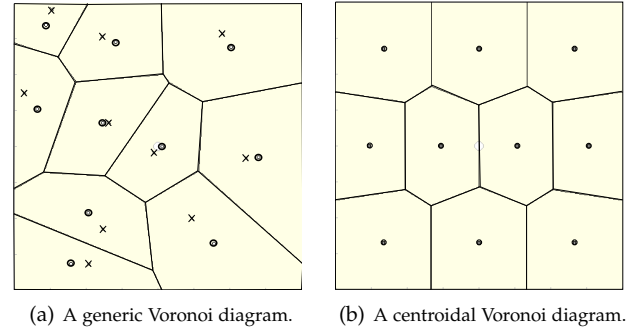


Fig. 4. The Voronoi tessellation of 10 points.

those Voronoi cells have different sizes. To achieve the load balance among those switches, we introduce the theory of Centroidal Voronoi Tessellation (CVT) [30] to further refine the coordinates of switches so that each Voronoi cell [29] has the equal size.

In Fig. 4(a) the crosses are the Voronoi sites and the circles are the centroids of the corresponding Voronoi regions. Note that the sites and the centroids do not coincide in Fig. 4(a). However, Fig. 4(b) shows a 10-point CVT, which can be viewed as an optimal partition corresponding to an optimal distribution of sites. That is, the circles are the sites for the Voronoi tessellation and the centroids of the corresponding Voronoi regions. Further, we hope that the coordinates of switches are also the centroids of the related Voronoi regions. After that, we can achieve the proper load balance when those data items are mapped into the virtual space evenly.

In geometry, a CVT [30] is a special type of Voronoi tessellation or Voronoi diagram, which is a partitioning of a plane into regions based on distance to points in a specific subset of the plane. The constraint for the CVT is simply that each Voronoi site must be the mass centroid for its corresponding Voronoi region. Given a region  $R$  and a density function  $\rho$ , defined in  $\Omega$ , the mass centroid  $r^*$  of the region  $R$  is defined by

$$r^* = \frac{\int_R r\rho(r)dr}{\int_R \rho(r)dr} \quad (1)$$

Given the number of sites  $n$ , a CVT is a minimizer (or a local minimizer) of the CVT *energy*, defined to be the square of the distance between each point in the region and its nearest site. Let  $\Omega$  be a metric space with distance function  $\phi$ . Assume that there are  $n$  sites, and  $(q_k)_{1 \leq k \leq n}$  be a site in the space  $\Omega$ . If  $\phi(r, P) = \inf\{\phi(r, q) | q \in P\}$  denotes the distance between the point  $r$  and the subset  $P$ , then we define a region  $R_k$  associated with the site  $q_k$  as follows.

$$R_k = \{r \in \Omega | \phi(r, q_k) \leq \phi(r, q_j), j = 1, \dots, n, j \neq k\} \quad (2)$$

That is, the region  $R_k$  is the set of all points in  $\Omega$  whose distance to  $q_k$  is not greater than their distance to the other sites  $q_j$ , where  $j$  is any index different from  $k$ . Accordingly, these regions are called Voronoi cells, and the diagram is a general Voronoi diagram. Furthermore, given a density

**Algorithm 1** *C-regulation*: refine the coordinates of switches in the virtual space while achieving the load balance.

**Require:** The coordinates of the switches  $Q$  achieved in Section 4.1.

**Ensure:** The updated coordinates of the switches  $Q^*$ .

- 1: Set  $Q^* \leftarrow Q$ ; set  $j_i=1$  for  $i=1, \dots, n$  where  $q_i \in Q$ ;
- 2: Obtain a random sample  $W$  of 1000 points from the virtual space  $\Omega$  that is constructed by the control plane with uniform probability;
- 3: For each point  $w \in W$ , find the  $q_i$  that is closest to  $w$ ; denote the index of that  $q_i$  by  $i^*$ ;
- 4: Set  $q_{i^*} \leftarrow \frac{j_{i^*} q_{i^*} + w}{j_{i^*} + 1}$  and  $j_{i^*} \leftarrow j_{i^*} + 1$ ; this new  $q_{i^*}$ , along with the unchanged  $q_i, i \neq i^*$ , form the new set of points  $Q^*$ . Note that  $j_i - 1$  equals the number of times that the point  $q_i$  has been updated.
- 5: If this new set of points meets some convergence criterion, terminate; otherwise, go back to step 2.

function  $\rho(\cdot)$  defined on  $\Omega$ , the formulation of the CVT energy is as follows:

$$F((q_i, R_i), i=1, \dots, n) = \sum_{i=1}^n \int_{r \in R_i} \rho(r) |r - q_i|^2 dr \quad (3)$$

Inspired by the CVT, we design the *C-regulation* method, as shown in Algorithm 1, to further refine the positions of switches obtained by the *Scoord* algorithm in Section 4.1. The *C-regulation* algorithm is a sampling technique, which supplies a discrete estimate of this CVT energy. Theorem 1 gives a necessary condition for the CVT energy  $F$  to be minimized, which means that the  $R_i$  is the Voronoi region corresponding to the switch's coordinate  $q_i$  and, simultaneously, the switch's coordinate  $q_i$  is the centroid of the corresponding Voronoi region  $R_i$ , for all  $1 \leq i \leq n$ . Based on Theorem 1, each time this *C-regulation* iteration is carried out, an attempt is made to modify the coordinates of switches in such a way that they are closer and closer to being the centroids of the Voronoi cells they generate. After that, the *C-regulation* algorithm can efficiently balance the Voronoi cell, and further achieve the load balance of switches, where the coordinates of switches are the sites of Voronoi cells.

The iteration will terminate when the CVT energy is lower than a given threshold. We set that the number of sample points is 1000 in each iteration, and that can be more. Note that the *C-regulation* method could require fewer iterations when more points are sampled in each iteration. However, more sample points will incur more computing time in each iteration. In addition, the number of iterations can also be set as the termination condition. During a iteration of the *C-regulation* algorithm, it should generally be the case that the CVT energy decreases from step to step. Furthermore, the impact of the number of iterations on the load balance is evaluated in Section 7.5.3. When the *C-regulation* algorithm terminates, we can achieve the updated coordinates of switches, which are indicated by the set of points  $Q^*$  in Algorithm 1.

**Theorem 1.** A necessary condition for the CVT energy  $F$  to be minimized is that the  $R_i$  is the Voronoi region corresponding to the switch's coordinate  $q_i$ , and the switch's coordinate  $q_i$  is the centroid of the corresponding Voronoi region  $R_i$ , for all  $1 \leq i \leq n$ .

*Proof:*

Given a region  $\Omega$ , a positive integer  $n$ , and a density function  $\rho(\cdot)$  defined on  $\Omega$ , let  $\{q_i\}_{i=1}^n$  denote the set of  $n$  switches' coordinates belonging to  $\Omega$  and let  $\{R_i\}_{i=1}^n$  denote any tessellation of  $\Omega$  into  $n$  regions. let

$$F((q_i, R_i), i=1, \dots, n) = \sum_{i=1}^n \int_{r \in R_i} \rho(r) |r - q_i|^2 dr \quad (4)$$

First, examine the first variation of  $F$  with respect to a single coordinate, say,  $q_j$ :

$$F(q_j + \epsilon \nu) - F(q_j) = \int_{r \in R_j} \rho(r) \{|r - q_j - \epsilon \nu|^2 - |r - q_j|^2\} dr \quad (5)$$

where we have not listed the fixed variables in the argument of  $F$  and where  $\nu$  is arbitrary such that  $q_j + \epsilon \nu \in \Omega$ . Then, when the CVT energy  $F$  is minimized, one easily finds that

$$q_j = \frac{\int_{R_j} r \rho(r) dr}{\int_{R_j} \rho(r) dr}, \quad (6)$$

by dividing Equation (5) by  $\epsilon$  and taking the limit as  $\epsilon \rightarrow 0$ .

Thus, according to the sense of Equation (1), the switch's coordinate  $q_j$  are the centroids of the regions  $R_j$  when the energy  $F$  is minimum. Next, let us hold the switches' coordinates  $\{q_i\}_{i=1}^n$  fixed and see what happens if we choose a tessellation  $\{\hat{R}_i\}_{i=1}^n$  other than the Voronoi tessellation  $\{\hat{R}_j\}_{j=1}^n$ . Let us compare the value of  $F((q_i, R_i), i=1, \dots, n)$  given by Equation (4) with that of

$$F((q_j, \hat{R}_j), j=1, \dots, n) = \sum_{j=1}^n \int_{r \in \hat{R}_j} \rho(r) |r - q_j|^2 dr \quad (7)$$

At a particular value of  $r$ ,

$$\rho(r) |r - q_j|^2 \leq \rho(r) |r - q_i|^2 \quad (8)$$

According to the sense of Equation (2), this result follows because  $r$  belongs to the Voronoi region  $\hat{R}_j$  corresponding to  $q_j$  and possibly not to the Voronoi region corresponding to  $q_i$ ; i.e.,  $r \in R_i$  but  $R_i$  is not necessarily the Voronoi region corresponding to  $q_i$ . Since  $\{R_i\}_{i=1}^n$  is not a Voronoi tessellation of  $\Omega$ , Equation (8) must hold with strict inequality over some measurable set of  $\Omega$ . Thus,

$$F((q_j, \hat{R}_j), j=1, \dots, n) < F((q_i, R_i), i=1, \dots, n) \quad (9)$$

so that  $F$  is minimized when the subsets  $R_i, i=1, \dots, k$ , are chosen to be the Voronoi regions associated with the switches' coordinates  $r_i, i=1, \dots, n$ .

Thus, Theorem 1 is proved.  $\square$

### 4.3 Multi-hop DT construction

To achieve the guaranteed delivery, the control plane constructs a multi-hop DT in the virtual space. As shown in Fig. 2(b), that is a multi-hop DT graph of 10 points. Furthermore, greedy routing in a multi-hop DT provides the property of guaranteed delivery [13], which is based on a rigorous theoretical foundation. For a given set of nodes in a 2D space, a triangulation is to construct edges between pairs of nodes such that the edges form a non-overlapping set of triangles that cover the convex hull of the nodes. Furthermore, a DT [31] in a 2D space is usually defined as a

triangulation such that the circumcircle of each triangle does not include any node other than the vertices of the triangle.

After obtaining the switches' positions in a set of points  $Q^*$ , a randomized incremental algorithm is designed to construct the DT  $DT(Q^*)$  in the 2D virtual space [32]. We first add an appropriate triangle boundingbox to contain  $P$ . The points in  $P$  are inserted in random order, and a DT corresponding to the current point set is maintained and updated throughout the whole process. Last, we remove the boundingbox and relative triangles which contains any vertex of the boundingbox triangle. Meanwhile, it is necessary to ensure that the union of all simplices in the triangulation is the convex hull of those points. Furthermore, greedy routing on a DT graph can achieve the guaranteed delivery [22]. That is, given a destination location  $p$ , the data packets always stop at a node that is closest to  $p$  among all nodes.

Considering the case of inserting  $v_i$ ,  $DT(v_1, v_2, \dots, v_{i-1})$  formed by inserting all previous points  $v_1, v_2, \dots, v_{i-1}$  is already a DT. The change caused by inserting  $v_i$  is adjusted and  $DT(v_1, v_2, \dots, v_{i-1}) \cup v_i$  is made a new  $DT(v_1, v_2, \dots, v_i)$ . The adjustment process is as follows. First, we determine which triangle (or edge)  $v_i$  falls on, and then connect  $v_i$  with the three vertices of the triangle to form three triangles (or connect the vertices of two triangles of the common edge to form four triangles). Since the newly generated edges and the original edges may not be Delaunay edges, a flipping [31] is conducted to make them all Delaunay edges to get  $DT(v_1, v_2, \dots, v_i)$ . Take  $DT(A, B, C, D)$  for example, we change the common edge  $\langle B, D \rangle$  to the common edge  $\langle A, C \rangle$  to produce two triangles that do meet the Delaunay condition when two original triangles do not meet the Delaunay condition [31]. This operation is called a flipping.

However, a key challenge is to ensure that each switch can transfer data packets to its DT neighbors note that a DT neighbor of a switch may not be the physical neighbor of the switch. Therefore, to achieve the guaranteed delivery, each switch maintains two kinds of flow entries in the GRED protocol, one makes it can forward requests to its physical neighbors, and another makes it forward requests to its multi-hop DT neighbors. Note that the switches that are not directly connected to some edge servers will not participate in the construction of the DT. Those switches are just used as the intermediate nodes to transfer data to the multi-hop DT neighbors. For a node  $u$ , each entry in its forwarding table  $F_u$  is a 4-tuple as follows.

$$\langle sour, pred, succ, dest \rangle,$$

which is a sequence of nodes with  $sour$  and  $dest$  being the source and destination nodes of a path, and  $pred$  and  $succ$  being node  $u$ 's predecessor and successor nodes in the path.  $F_u$  is used to forward packets to multi-hop DT neighbors. For a specific tuple  $t$ , we use  $t.sour$ ,  $t.pred$ ,  $t.succ$ , and  $t.dest$  to denote the corresponding nodes in the tuple  $t$ . Although greedy routing does not always find a shortest route, the quality of the greedy route is often very good. The length of an optimal route between a pair of nodes on a DT is within a constant time of the direct distance [33].

---

### Algorithm 2 GRED( $u, d$ ) forwarding at switch $u$ .

---

- 1: For each physical neighbor  $v$ ,  $R_v \leftarrow ED(v, d)$ , Euclidean distance between  $v$  and  $d$ ;
  - 2: For each DT neighbor  $\tilde{v}$ ,  $R_{\tilde{v}} \leftarrow ED(\tilde{v}, d)$ ;
  - 3: Let  $v^*$  be the neighbor where  $R_{v^*} = \min\{R_v, R_{\tilde{v}}\}$ ;
  - 4: **if**  $R_{v^*} < ED(u, d)$  **then**
  - 5:     Send the packet to  $v^*$  directly or by the multi-hop path;
  - 6: **else**
  - 7:     Switch  $u$  is closest to  $d$ , and determines a unique edge server to place the data;
  - 8: **end if**
- 

## 5 DATA PLACEMENT AND RETRIEVAL USING GRED

In this section, we detail how the GRED is designed to support data placement and retrieval services. The GRED protocol would first forward data packets to the switch, which is closest to the data position in the virtual space. Then, the switch would determine a unique edge server to place the data. Furthermore, we introduce how to use the GRED protocol to retrieve a data item.

### 5.1 Placing data in the edge network

In GRED, the switches are associated with their coordinates in the virtual space, which is maintained by the control plane. A switch knows its own coordinates, its physical neighbors' coordinates, and the coordinates of its DT neighbors. The Euclidean distance between any two switches can be calculated from their coordinates where the network-wide distance has been embedded in Section 4.1. The key idea of GRED forwarding at a switch, say  $u$ , is conceptually simple: For a data with ID  $d$ , the place to store the data is position  $H(d)$ , which will be converted to the coordinate in the virtual space, as shown in Section 3.  $u$  forwards the packet to the DT-neighbor switch closest to  $H(d)$ . If the neighbor is a physical neighbor, the packet is directly forwarded; else, the packet is forwarded via a virtual link, to a DT neighbor closest to  $H(d)$ . If there is no neighbor of  $u$  closer to  $H(d)$  than  $u$  itself, it is proved that  $u$  is the switch closest to  $H(d)$  [33]. When the data arrives at the switch closest to  $H(d)$ , the switch determines a unique edge server to place the data. The detailed algorithm is presented in Algorithm 2.

**Transfer in a virtual link.** Consider a switch  $u$  that has received a data  $d$  to forward. Switch  $u$  stores it with the format:  $d = \langle d.dest, d.sour, d.relay, d.data \rangle$  in a local data structure. When  $d.relay \neq null$ , data  $d$  is traversing a virtual link. Note that  $d.dest$  is the end switch of the virtual link,  $d.sour$  is the source switch,  $d.relay$  is the relay switch, and  $d.data$  is the payload of the data. When switch  $u$  receives a packet that is being forwarded in a virtual link, the packet is processed as follows. When  $u = d.dest$ , switch  $u$  is the endpoint of the virtual link, and continues to forward the data based on Algorithm 2. When  $u = d.relay$ , switch  $u$  first finds tuple  $t$  from the forwarding table  $F_u$  with  $t.dest = d.dest$  where  $F_u$  is defined in Section 4.3. Then, switch  $u$  revises  $d.relay = t.succ$  based on the matched tuple  $t$ . The last step in switch  $u$  is to transmit the data to  $d.relay$ . Based on this setting, messages can be forwarded to a DT neighbor of a switch. However, it is worth noting that a

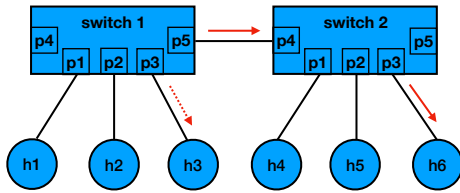


Fig. 5. Data item that should be placed in server  $h3$  is placed in server  $h6$  when server  $h3$  would be overloaded.

global minimum may not be unique, for those data mapped to a Voronoi edge in Fig. 4(b). The tie can be broken by ranking the  $x$  coordinate, then  $y$  coordinate.

## 5.2 Determining the placement server

Based on the above analysis, GRED can ensure that a data item can be forwarded to a unique switch, whose position is closest to the position of the data. Furthermore, the switch determines a unique server to place the data. Assume that switch  $u$  is closest to the data position in the virtual space, and switch  $u$  is directly interconnected with  $s$  edge servers. In GRED protocol, switch  $u$  maintains a serial number for each edge server from 0 to  $s-1$ . Then, switch  $u$  transmits the data with the identifier  $d$  to the server whose serial number is  $[H(d) \bmod s]$  where we still use a uniform hash function [23]. Furthermore, the method to determine the server can efficient balance the load among those edge servers since the hash function can map the expected inputs as evenly as possible over its output range.

**The range extension.** Consider that edge servers could be heterogeneous. Some edge servers with low storage capacity would be overloaded when switches connect to different numbers of edge servers with heterogeneous capacity. To solve this problem, we further extend the management range of the switches. The management range of a switch is determined by the edge servers that the switch can place data. In prior discussion, the management range is one-hop. That is, the data whose position is closest to a switch position would be placed in the edge server directly connected to the corresponding switch. Furthermore, GRED allows that a switch can manage servers with more than one hop. When the upper layer application finds that an edge server would be overloaded, the corresponding switch sends an extending request to the control plane, which can be achieved in the context of SDN [8]. Accordingly, the control plane assigns the edge server with the most

remaining capacity from the physical neighbor switches to take over the corresponding storage load. To enable this, the control plane needs to update the corresponding forwarding entries into the related switches.

As shown in Fig. 5, when the server  $h3$  that connected to switch 1 would be overloaded, the switch 1 sends an extending request to the control plane. Then, the control plane assigns server  $h6$  to take over the load of server  $h3$  where the edge server  $h6$  is connected to switch 2. Before that, for switch 1, the data that should be placed in server  $h3$  would be forwarded to port  $p3$  based on the flow entry in Table 1. However, the data is forwarded to port  $p5$  after that the control plane replaces the forwarding entry in Table 1 with the flow entry in Table 2. Table 2 shows that switch 1 first sets the destination address of the data as the address of server  $h6$ , and then forwards the data to port  $p5$ , when the destination address of the data is the address of server  $h3$ . Meanwhile, switch 2 also receives the corresponding forwarding entry, which indicates to forward the related data to its edge server  $h6$ . Furthermore, when some edge servers in switch 2's range would be overloaded, switch 2 will also send an extending request to the control plane. Therefore, the range extension can efficiently avoid the overload of edge servers and share the resources of multiple edge servers.

In addition, consider that the data placement in edge servers is not everlasting. That is, the overloaded edge server could become underloaded again since some data could be invalid or migrated to the Cloud. In this case, the edge server will first retrieve the data, which should be placed in the edge server, but is placed in other edge servers. When all the corresponding data has been retrieved, the corresponding extended forwarding entries will also be deleted from the related switches.

## 5.3 Data retrieval using GRED

So far, we have introduced the procedure of data placement. The data retrieval using GRED is similar to the data placement. The retrieval is also to use the data identifier, and each switch greedily forwards the retrieval request to the switch whose position is closest to the data position in the virtual space. Furthermore, the switch uses the same method shown in Section 5.2 to determine the edge server for responding to the retrieval request. However, the key challenge is how to determine the edge server that has stored a data when the corresponding switch has extended its management range.

As shown in Fig. 5, the data that should be placed in server  $h3$  that is connected to switch 1 is forwarded to server  $h6$  connected to switch 2 when switch 1 extends its management range. In this case, when we retrieve a data that is directed to the edge server  $h3$  based on the value of  $[H(d) \bmod s]$ , we could not determine that the data has been placed in server  $h3$  connected to switch 1 or server  $h6$  connected to switch 2. Therefore, to efficiently retrieve a data, the retrieval request is forwarded to the two edge servers at the same time, and the edge server that has stored the data will respond to the retrieval request. Note that a tag is used in the packet header to indicate a placement/retrieval request. After that, we can ensure to

TABLE 1  
The flow entry in switch 1 before updating.

	Match	Action
1	$d.dest=h3.address$	Output: port $p3$

TABLE 2  
The flow entry in switch 1 after updating.

	Match	Action
1	$d.dest=h3.address$	Set: $d.dest=h6.address$ ; Output: port $p5$



efficiently locate a data that has been placed in the edge network when a data retrieval request is received.

## 6 DISCUSSION

**The network dynamic.** Consider that some edge nodes could be added into the edge network. Meanwhile, some failures of switches or edge nodes could result in that some edge nodes leave from the edge network. Therefore, the GRED is required to accommodate the network dynamic. Recall that we utilize an incremental method to construct the DT graph in the control plane in Section 4.3. When an edge node is added in the edge network, some edges will be added in the DT graph to connect the new edge node and its neighbors, which have existed in the DT graph. It is worth noting that the new edge node has no effect on the other edge nodes. It only affects its neighbors. First, the control plane will add the corresponding forwarding entries into the new edge nodes and its neighbors. Then, those data in the neighboring edge nodes of the new edge node will be calculated again. If those data is closest to the new edge node, they will be forwarded to the new edge node. Furthermore, when an edge node leaves from the edge network, the related edges between it and its neighbors will be deleted, and then some new edges will be added between those neighbors to form a new DT graph. After that, those related data will be stored in those neighbors based on their positions in the virtual space, which has been described in Section 5.1.

**Data copies.** The data copies are fundamental for the fault tolerance. Meanwhile, multiple data copies can also help to achieve better performance. Therefore, it is necessary for the GRED to support multiple data copies in the edge network. Recall that we store a data item in the edge network by hashing its ID. Furthermore, when there exists multiple data copies, it is required to add a serial number for each data copy. Then, the ID and the serial number are concatenated to form a new string. By hashing the new string, we can achieve the position of the corresponding data copy in the virtual space. After that, the data copy can be stored in the edge network based on the scheme in Section 5.1. An advantage of the GRED is that it is easy to determine which copy is closest to the access point. Consider that we have embedded the network-wide distance between switches into the Euclidean distance between the related two points in the virtual space in Section 4.1. Therefore, we can know which copy is closest to the access point by calculating their distances to the access point in the virtual space after embedding the network distance in Section 4.1.

## 7 PERFORMANCE EVALUATION

In this section, we first introduce the implementation and evaluation of the GRED on a small-size testbed. Then, we conduct large-scale simulations to evaluate the performance of the GRED.

### 7.1 Implementation using P4

We have implemented a prototype of GRED, including all switch data plane and control plane features described in Section 3, where the switch data plane is written in P4 [11],

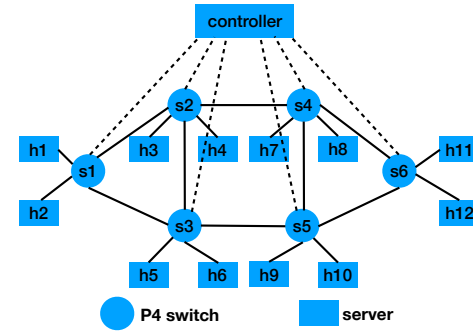


Fig. 6. Prototype with 1 controller, 6 P4 switches and 12 servers.

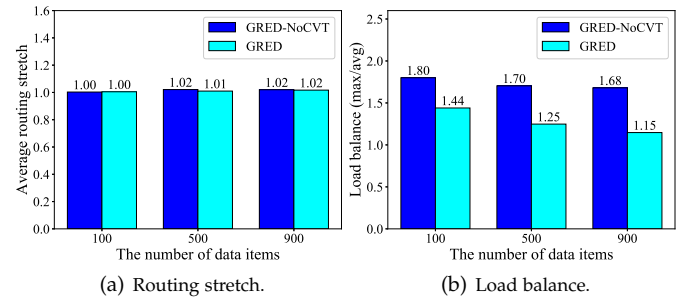


Fig. 7. The performance of the GRED protocol under different settings.

and the function in the control plane is written in Java. The P4 compiler generates Thrift APIs for the controller to insert the forwarding entries into the switches. The P4 switch supports a programmable parser to allow new headers to be defined. Meanwhile, multiple match+action stages [11] are designed in series to achieve the neighboring switch whose position is closest to the position of the data. The P4 switch calculates the distance from a neighbor to the data in the virtual space in a match+action stage. The topology of our prototype is shown in Fig. 6. Our testbed consists of 1 controller and 6 P4 switches, where each switch connects to 2 servers. We use those servers to generate data requests including the data placement/retrieval requests. Furthermore, we evaluate the routing stretch and the load balance of the GRED protocol on our prototype. We implemented two variants of the GRED protocol including the GRED-NoCVT protocol and the GRED protocol on our testbed. The GRED protocol sets the number of iterations is 50 for the C-regulation method, which is introduced in Section 4.2. GRED-NoCVT indicates the positions of switches are only generated by the *Scoord* algorithm in Section 4.1, and not refined by the C-regulation method.

We first evaluate the performance of the GRED protocol based on our testbed. Fig 7(a) shows that the average routing stretches of GRED-NoCVT and GRED are close to 1, which is the optimal value of the routing stretch. However, Fig 7(b) shows that GRED achieves significantly better load balance than GRED-NoCVT due to the lower *max/avg* value, which is used to quantify the load balance of a networked storage system. The value of *max* is the number of data items received by the most loaded edge server, and the value of *avg* means the average load of all edge servers. The optimal value of *max/avg* is 1, which indicates perfect load balancing.

Furthermore, we test the average response delay of the GRED protocol where we have placed some data items

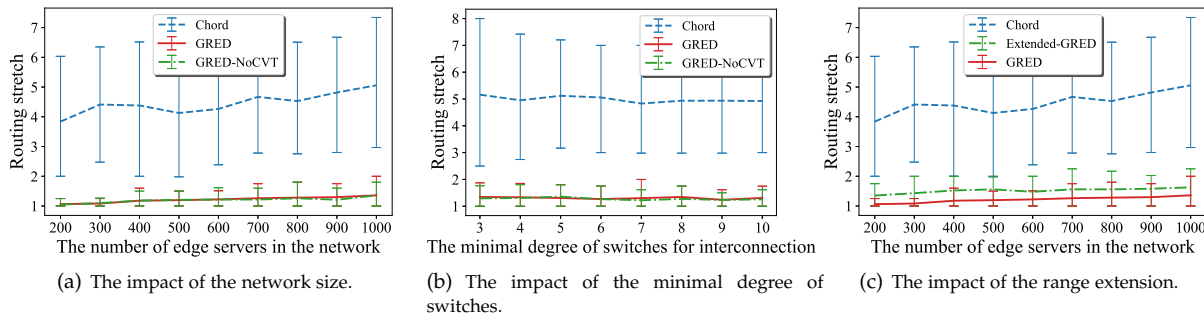


Fig. 8. Routing stretch comparison under different schemes.

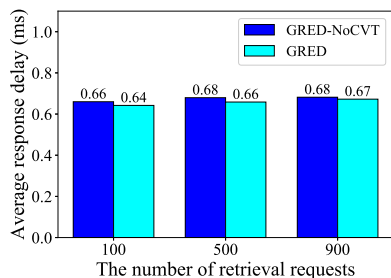


Fig. 9. The response delay under different number of retrieval requests.

in our testbed and then generated some data retrieval requests. Fig 9 shows that the average response delay of those retrieval requests. We can find that the average response delays of the two GRED variants are similar, and the average response delay has a modest change when we send the different number of retrieval requests. The routing stretch would affect the average response delay. Recall that the two GRED variants all have low routing stretches in Fig 7(a). Therefore, we can find that the response delay is low in Fig 9. That is, the GRED protocol can quickly respond to those retrieval requests in edge computing. However, it is worth noting that the network size is small since our testbed just consists of 6 P4 switches and 12 edge servers. So, we further conduct massive simulations to evaluate the performance of the GRED protocol including the routing stretch and the load balance in the next section.

## 7.2 The setting of large-scale simulations

In simulations, unless otherwise specified, we use BRITE [34] with the Waxman model to generate synthetic topologies at the switch level where each switch connects to 10 edge servers. Switches could connect to different numbers of edge servers or servers with different capacity. Then, we compare the GRED protocol with the Chord [12] protocol, which can locate data in a peer-to-peer network. The GRED protocol includes two variants: GRED and GRED-NoCVT (without CVT). We use two performance metrics to evaluate the performance of GRED as follows.

- **Routing stretch.** The routing stretch value is defined to be the ratio of the hop count in the selected route to the hop count in the shortest route between a pair of source and destination nodes.
- **Load balance.** The  $max/avg$  metric quantifies the load balance, defined as the ratio of the number of data items received by the most loaded edge server ( $max$ ) to the average load of all edge servers ( $avg$ ).

We evaluate the routing stretch of GRED by varying the number of switches and the minimal degree of switches for interconnection. In each setting of the network, we randomly generate 100 data items to be placed in the network and randomly select an access point for each data. Each point in Fig. 8 is the average of 100 routing stretches where each error bar is constructed using a 90% confidence interval of the mean. Furthermore, we evaluate the load balance of GRED varying the number of switches and the amount of data. Meanwhile, we evaluate the impact of the number of iterations of the C-regulation method on the load balance of GRED.

## 7.3 Routing stretch

### 7.3.1 Varying network size

We first evaluate the impact of the network size on the routing stretch. Fig. 8(a) shows the routing stretches of Chord, GRED, and GRED-NoCVT. In Fig. 8(a), GRED and GRED-NoCVT achieve significantly lower routing stretches than Chord. It is because that the Chord takes  $O(\log n)$  overlay hops to retrieve the data while the GRED costs only one overlay hop to get the data. The average routing stretch of Chord is higher than 3.5 under any network size in our experiments. However, the average routing stretches of GRED and GRED-NoCVT are all lower than 1.5. It means that GRED uses <30% routing path lengths compared to using Chord. It is worth noting that shorter routing path indicates less bandwidth consumption and lower latency to place/retrieve data. Meanwhile, we can see that GRED has a little higher routing stretch than GRED-NoCVT in some cases. It is because the C-regulation method has influence on the distances between switches, which can be preserved as well as possible after using the *Scoord* algorithm in Section 4.1.

### 7.3.2 Varying the minimum degree of switches

We evaluate the impact of the minimal degree of switches for interconnection on the routing stretch. The network employs 100 switches and 1000 edge servers, and the minimal degree of switches for interconnection varies from 3 to 10. Fig. 8(b) shows that GRED and GRED-NoCVT achieve obviously lower routing stretches than the Chord protocol. In Fig. 8(b), we can see that the degree of switches for interconnection has a modest impact on the routing stretch for the same protocol. Meanwhile, Fig. 8(b) shows that the routing stretch slightly decreases as the increase of the minimal degree of switches. When the switches provide

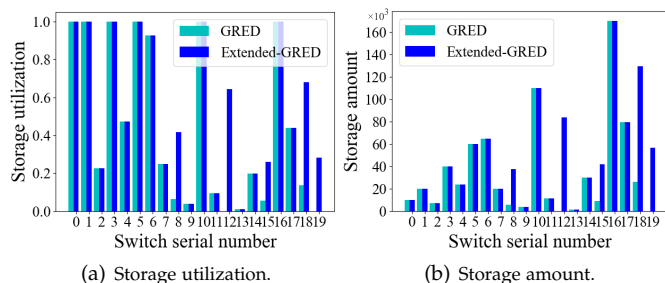


Fig. 10. The performance of the GRED protocol with the range extension.

more ports for interconnection, greedy routing has a higher possibility to find the shortest path.

### 7.3.3 Range extension

When an edge server will be overloaded, the corresponding switch needs to extend its management range. That is, the switch forwards data to the edge server connected to the neighboring switch. Range extension may increase the routing stretch. We compare the routing stretch achieved by GRED and the extended-GRED protocol where the number of iterations is 50 for the C-regulation method. The extended-GRED denotes the data would be placed in the edge server connected to the neighbor switch of the destination switch. We placed 100 data items to achieve the average routing stretch under each setting of the network size. Fig. 8(c) shows that the extended-GRED protocol achieves slightly higher routing stretch than GRED. However, the routing stretch of the extended-GRED is still significantly lower than Chord.

Furthermore, we evaluate how the range extension improve the storage utilization of edge servers. 1 million data items are stored in the edge computing environment where 20 switches exist in the network, and the storage capacities of edge servers connected with switches vary from 10K to 200K. Note that the edge servers will drop the received data when it is overloaded under the GRED without the range extension. From Fig. 10(a), we can see that the storage utilizations of edge servers connected with switches {8, 12, 15, 18, 19} are significantly improved. The storage utilization in switch 12 increases 60% more under the extended-GRED compared to the GRED in Fig. 10(a). Because the switch will forward the data items to its neighbors when those edge servers connected to the corresponding switch are overloaded under the extended-GRED. Fig. 10(b) shows that the corresponding edge servers store more data items under the extended-GRED than that of the GRED. Meanwhile, we can find that the improvement to the edge servers with full storage is small. Therefore, the extended-GRED can efficiently utilize the left storage resources in some edge servers, and further improve the service capacity of edge infrastructure for cooperative edge clouds.

## 7.4 The number of forwarding table entries

In this section, we show the number of forwarding table entries per switch for the GRED protocol under different network sizes. In Figure 11, each point indicates the average number of forwarding table entries over all switches, where the error bar is constructed using a 90% confidence

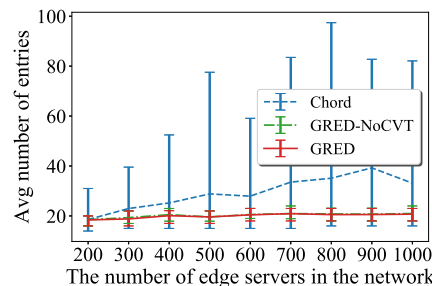


Fig. 11. The number of forwarding table entries under different schemes.

interval of the mean. Figure 11 shows that the number of forwarding table entries goes up as the increase of the network size under the Chord protocol. However, we can see that the increase of the average number of forwarding entries is modest as the increase of the network size under the GRED and GRED-NoCVT protocols from Figure 11. That is, the GRED protocol only needs a few forwarding entries to achieve the data placement and retrieval services. In addition, under the Chord protocol, not only switches employ forwarding entries to support the data placement and retrieval service, but also edge servers need to maintain finger tables to achieve data location. Therefore, GRED has the obvious advantage in scalability since the number of forwarding table entries is independent of the network size and the number of flows in the edge network.

## 7.5 Load balance

### 7.5.1 Varying the network Size

We first evaluate the impact of the network size on the load balance under different protocols where the number of edge servers varies from 200 to 1000. Fig. 12(a) shows that GRED ( $T=10$ ) and GRED ( $T=50$ ),  $T$  is the number of iterations, achieve significantly better load balance than Chord due to the lower value of  $max/avg$ . In Fig. 12(a), the value of  $max/avg$  goes up as the increase of the network size in Chord. However, we observe very little increase for GRED ( $T=10$ ) and GRED ( $T=50$ ) in Fig. 12(a). Fig. 12(a) shows that GRED ( $T=50$ ) achieves better load balance than GRED ( $T=10$ ), which means that the GRED protocol can achieve better load balance by increasing the number of iterations.

### 7.5.2 Varying the amount of data

We vary the amount of the placed data from 100,000 to 1,000,000 where 1000 edge servers are deployed in the network. Fig. 12(b) shows that GRED ( $T=50$ ) achieves the best load balance among the three protocols. We can see that the Chord protocol has the worst load balance because the value of  $max/avg$  is higher than 6. Meanwhile, we can also see that the value of  $max/avg$  for GRED ( $T=10$ ) is lower than 2.5, and further the value of GRED ( $T=50$ ) is lower than 2. Note that the value of  $max/avg$  is lower and better, and the optimal value for load balance is 1. Therefore, the GRED protocol can achieve the proper load balance among edge servers.

### 7.5.3 Varying the number of iterations

In this section, we test the impact of the number of iterations  $T$  on the load balance. Note that the number of iterations



Fig. 12. Comparison of load balance under different schemes.

$T$  for the C-regulation method will affect the positions of switches in the virtual space, and further affect the load balance of the GRED protocol. The setting of the network is the same as the setting in Section 7.5.2, and we placed 100,000 data items in the network. Note that the Chord and the GRED-NoCVT are independent of  $T$ . Therefore, Fig. 12(c) shows that  $T$  has no influence on Chord and GRED-NoCVT. Furthermore, we can see that the value of  $max/avg$  decreases as the increase of  $T$  for the GRED protocol in Fig.12(c). That means that the GRED protocol can achieve better load balance when  $T$  increases. Meanwhile, Fig. 12(c) shows that GRED-NoCVT can also achieve better load balance than Chord even if GRED-NoCVT did not use the C-regulation method to refine the positions of switches. Furthermore, we can see that the value of  $max/avg$  is lower than 2 when  $T$  is more than 20 in Fig. 12(c). We also find that the value of  $max/avg$  stops to decrease when  $T$  is more than 70 in Fig. 12(c). It means that the C-regulation method has found the optimal positions of switches in the virtual space to achieve the proper load balance when  $T=70$ . After that, the increase of  $T$  has little improvement on the load balance of GRED.

## 8 RELATED WORK

In recent years, a new trend in computing is happening with the function of Clouds being increasingly moving towards the network edges [35]. It is estimated that tens of billions of edge devices will be deployed in the near future [36]. The computing and storage resources are placed at the edge of the Internet to provide low-latency services for those edge devices. Zeng et al. study how to effectively and economically utilize the idle resources in volunteer vehicles to handle the overloaded tasks in VEC servers [37]. Chen et al. propose a light-weight radio frequency fingerprinting identification (RFFID) scheme to realize authentications for a large number of resource-constrained terminals under the mobile edge computing (MEC) scenario without relying on encryption-based methods [38]. Liao et al. investigate the security threats in mobile edge computing (MEC) of Internet of things, and propose a deep-learning (DL)-based physical (PHY) layer authentication scheme [39]. Wu et al. discuss the roles and opportunities that information and communications technologies play in pursuing the sustainable development goals [40]. Atat et al. present the cyber-physical systems taxonomy via providing a broad overview of data collection, storage, access, processing, and analysis [41].

In edge computing, edge servers perform computing offloading, data storage, caching and processing, as well as

distribute request and delivery service [1]. Yang et al. propose a data centric design where data become self-sufficient entities that are stored, referenced independently from their producers [2]. However, a number of challenges need to be addressed in edge computing. First, Mobility is an intrinsic trait of many mobile applications. In those applications, the edge servers could exploit the movement and trajectory of edge users to improve the efficiency of handling users' computation requests. Some mobility models were proposed [42], which characterize the mobility by a sequence of networks that users can connect to and a two-dimensional location-time workflow, respectively. In addition, mobility management for edge computing was integrated with traffic control in [43] to provide better experience for users. Note that most of the existing works focused on optimizing mobility-aware edge server selection. However, to achieve better user experience and higher network-wide profit, we propose the GRED protocol to efficient locate data for edge users wherever users access the network.

On the other hand, edge servers with limited computational resources may be overloaded when they have to serve a large number of edge users. In such cases, The burdens on an edge server can be lightened via peer-to-peer cooperative edge servers [36]. Xia et al. investigate the collaborative caching problem in the EC environment with the aim to minimize the system cost including data caching cost, data migration cost, and quality-of-service (QoS) penalty [44]. Gharaibeh et al. study the collaborative caching problem for a multicell-coordinated system from the point of view of minimizing the total cost paid by the content providers [45]. Resource sharing via the cooperation of edge servers can not only improve the resource utilization, but also provide more resources for edge users to enhance their user experience. The resource sharing framework was originally proposed in reference [46], which includes components such as resource allocation, revenue management and service provider cooperation. The framework was extended in [47], which considered both the local and remote resource sharing. Server cooperation can significantly improve the computation efficiency and resource utilization at edge servers. More importantly, it can balance the computation and storage load distribution over the networks so as to reduce sum response latency. Therefore, in this paper, we propose the GRED protocol to efficient distribute and locate the data over the cooperative edge clouds.

## 9 CONCLUSION

Mobile edge computing needs to provide the data placement and retrieval services for many emerging applications

such as IoT. However, it remains an open problem. A key challenge to enable this is to efficiently locate the data in the edge network. GRED solves this challenging problem by offering a powerful primitive: given a data identifier, it determines the edge server responsible for the data placement and retrieval, and does so efficiently. Attractive features of GRED include its routing simplicity, provable correctness, low routing stretch, and proper load balance. Our theoretical analysis, simulations, and experimental results confirm that the effectiveness and efficiency of GRED. GRED uses <30% routing cost and achieves better load balance among edge clouds compared to using Chord, a well-known DHT. We believe that GRED will be a valuable component for mobile edge computing considering the user mobility and the cooperation among edge clouds.

## ACKNOWLEDGMENTS

This work is partially supported by the National Key Research and Development Program of China under Grant No. 2018YFE0207600, the National Natural Science Foundation of China under Grant No. U19B2024 and 62002284, and the Tianjin Science and Technology Foundation under Grant No. 18ZXJMTG00290.

## REFERENCES

- [1] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge computing: Vision and challenges," *IEEE Internet of Things Journal*, vol. 3, no. 5, pp. 637–646, 2016.
- [2] Y. Yang, "A vision towards pervasive edge computing," in *Proceedings of the 22nd International ACM Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems*, ser. MSWIM 19, 2019, p. 1.
- [3] M. Chiang and T. Zhang, "Fog and IoT: An Overview of Research Opportunities," *IEEE Internet of Things Journal*, vol. 3, no. 6, pp. 854–864, 2016.
- [4] Y. Elkhatib, B. Porter, H. B. Ribeiro, M. F. Zhani, J. Qadir, and E. Rivire, "On using micro-clouds to deliver the fog," *IEEE Internet Computing*, vol. 21, no. 2, pp. 8–15, 2017.
- [5] Y. Zeng, Y. Huang, Z. Liu, and Y. Yang, "Joint online edge caching and load balancing for mobile data offloading in 5g networks," in *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*, 2019, pp. 923–933.
- [6] F. Bonomi, R. Milito, J. Zhu, and S. Addepalli, "Fog computing and its role in the internet of things," in *Proc. of the 1st MCC Workshop on Mobile Cloud Computing*, August 2012.
- [7] J. Xie, C. Qian, D. Guo, X. Li, S. Shi, and H. Chen, "Efficient data placement and retrieval services in edge computing," in *Proc. of IEEE ICDCS*, July 2019.
- [8] D. Kreutz, F. M. V. Ramos, P. E. Verissimo, C. E. Rothenberg, S. Azodolmolky, and S. Uhlig, "Software-defined networking: A comprehensive survey," *Proceedings of the IEEE*, vol. 103, no. 1, pp. 14–76, 2015.
- [9] Y. Zeng, S. Guo, G. Liu, P. Li, and Y. Yang, "Energy-efficient device activation, rule installation and data transmission in software defined dns," *IEEE Transactions on Cloud Computing*, pp. 1–1, 2019.
- [10] J. Xie, D. Guo, Z. Hu, T. Qu, and P. Lv, "Control plane of software defined networks: A survey," *Computer Communications*, vol. 67, pp. 1–10, Aug. 2015.
- [11] P. Bosshart, D. Daly, G. Gibb, M. Izzard, N. McKeown, J. Rexford, C. Schlesinger, D. Talayco, A. Vahdat, G. Varghese, and D. Walker, "P4: Programming protocol-independent packet processors," *SIGCOMM Comput. Commun. Rev.*, vol. 44, no. 3, pp. 87–95, Jul. 2014.
- [12] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan, "Chord: A scalable peer-to-peer lookup service for internet applications," in *Proceedings of ACM SIGCOMM*, 2001, pp. 149–160.
- [13] S. S. Lam and C. Qian, "Geographic routing in d-dimensional spaces with guaranteed delivery and low stretch," *IEEE/ACM Trans. Netw.*, vol. 21, no. 2, pp. 663–677, 2013.
- [14] P. Mach and Z. Becvar, "Mobile edge computing: A survey on architecture and computation offloading," *IEEE Communications Surveys Tutorials*, vol. 19, no. 3, pp. 1628–1656, 2017.
- [15] Y. Huang, J. Zhang, J. Duan, B. Xiao, F. Ye, and Y. Yang, "Resource allocation and consensus on edge blockchain in pervasive edge computing environments," in *2019 IEEE 39th International Conference on Distributed Computing Systems (ICDCS)*, 2019, pp. 1476–1486.
- [16] R. Cox, A. Muthitacharoen, and R. T. Morris, "Serving dns using a peer-to-peer lookup service," in *International Workshop on Peer-To-Peer Systems*. Springer, 2002, pp. 155–165.
- [17] E. K. Lua, J. Crowcroft, M. Pias, R. Sharma, S. Lim *et al.*, "A survey and comparison of peer-to-peer overlay network schemes," *IEEE Communications Surveys and Tutorials*, vol. 7, no. 1-4, pp. 72–93, 2005.
- [18] Y. Liu, X. Liu, L. Xiao, L. M. Ni, and X. Zhang, "Location-aware topology matching in P2P systems," in *Proc. of IEEE INFOCOM*, 2004.
- [19] M. Caesar, M. Castro, E. B. Nightingale, G. O'Shea, and A. Rowstron, "Virtual Ring Routing: Networking Routing Inspired by DHTs," in *Proceedings of ACM Sigcomm*, 2006.
- [20] A. T. Mizrak, Y. Cheng, V. Kumar, and S. Savage, "Structured superpeers: Leveraging heterogeneity to provide constant-time lookup," in *Proc. of the Third IEEE WIAPP*, 2003, pp. 104–111.
- [21] V. Ramasubramanian and E. G. Sirer, "Beehive: O (1) lookup performance for power-law query distributions in peer-to-peer overlays," in *Proc. of USENIX NSDI*, 2004, pp. 99–112.
- [22] C. Qian and S. S. Lam, "ROME: Routing On Metropolitan-scale Ethernet," in *Proc. of IEEE ICNP*, 2012, pp. 1–10.
- [23] A. Biryukov, M. Lamberger, F. Mendel, and I. Nikolić, "Second-order differential collisions for reduced sha-256," in *Advances in Cryptology – ASIACRYPT*, 2011, pp. 270–287.
- [24] P. Berde, M. Gerola, J. Hart, Y. Higuchi, M. Kobayashi, T. Koide, B. Lantz, B. O'Connor, P. Radoslavov, W. Snow, and G. Parulkar, "Onos: Towards an open, distributed sdn os," in *Proceedings of the Third Workshop on Hot Topics in Software Defined Networking*, ser. HotSDN, 2014, pp. 1–6.
- [25] C. Qian and S. S. Lam, "Greedy routing by network distance embedding," *IEEE/ACM Transactions on Networking*, vol. 24, no. 4, pp. 2100–2113, 2016.
- [26] I. Borg and P. J. Groenen, *Modern multidimensional scaling: Theory and applications*. Springer Science & Business Media, 2005.
- [27] J. Xie, C. Qian, D. Guo, M. Wang, S. Shi, and H. Chen, "Efficient indexing mechanism for unstructured data sharing systems in edge computing," in *Proc. of IEEE INFOCOM*, April 2019, pp. 1–9.
- [28] F. Wickelmaier, "An introduction to mds," Sound Quality Research Unit, Aalborg University, Denmark, Tech. Rep., 2003.
- [29] S. Fortune, "Voronoi diagrams and Delaunay triangulations," in *Handbook of Discrete and Computational Geometry*, 2nd ed., J. E. Goodman and J. O'Rourke, Eds. CRC Press, 2004.
- [30] Q. Du, V. Faber, and M. Gunzburger, "Centroidal voronoi tessellations: Applications and algorithms," *SIAM Review*, vol. 41, no. 4, pp. 637–676, 1999.
- [31] J. A. De Loera, J. Rambau, and F. Santos, *Triangulations Structures for algorithms and applications*. Springer, 2010.
- [32] L. J. Guibas, D. E. Knuth, and M. Sharir, "Randomized incremental construction of delaunay and voronoi diagrams," *Algorithmica*, vol. 7, no. 1, pp. 381–413, 1992.
- [33] D. Y. Lee and S. S. Lam, "Protocol design for dynamic delaunay triangulation," in *Proc. of 27th IEEE ICDCS*, June 2007.
- [34] A. Medina, A. Lakhina, I. Matta, and J. Byers, "Brite: An approach to universal topology generation," in *Proc. 9th International Symposium on MASCOTS*, Cincinnati, OH, USA, August 2001.
- [35] K. Poularakis, J. Llorca, A. M. Tulino, I. Taylor, and L. Tassiulas, "Joint service placement and request routing in multi-cell mobile edge computing networks," in *Proc. of IEEE INFOCOM*, April 2019, pp. 10–18.
- [36] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," *IEEE Communications Surveys Tutorials*, vol. 19, no. 4, pp. 2322–2358, 2017.
- [37] F. Zeng, Q. Chen, L. Meng, and J. Wu, "Volunteer assisted collaborative offloading and resource allocation in vehicular edge computing," *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–11, 2020.
- [38] S. Chen, H. Wen, J. Wu, A. Xu, Y. Jiang, H. Song, and Y. Chen, "Radio frequency fingerprint-based intelligent mobile edge com-

puting for internet of things authentication," *Sensors*, vol. 19, no. 16, p. 3610, 2019.

- [39] R. Liao, H. Wen, J. Wu, F. Pan, A. Xu, H. Song, F. Xie, Y. Jiang, and M. Cao, "Security enhancement for mobile edge computing through physical layer authentication," *IEEE Access*, vol. 7, pp. 116 390–116 401, 2019.
- [40] J. Wu, S. Guo, H. Huang, W. Liu, and Y. Xiang, "Information and communications technologies for sustainable development goals: State-of-the-art, needs and perspectives," *IEEE Communications Surveys & Tutorials*, vol. 20, no. 3, pp. 2389–2406, 2018.
- [41] R. Atat, L. Liu, J. Wu, G. Li, C. Ye, and Y. Yang, "Big data meet cyber-physical systems: A panoramic survey," *IEEE Access*, vol. 6, pp. 73 603–73 636, 2018.
- [42] K. Lee and I. Shin, "User mobility model based computation offloading decision for mobile cloud," *Journal of Computing Science and Engineering*, vol. 9, no. 3, pp. 155–162, 2015.
- [43] A. Prasad, P. Lundén, M. Moisisio, M. A. Uusitalo, and Z. Li, "Efficient mobility and traffic management for delay tolerant cloud data in 5g networks," in *Personal, Indoor, and Mobile Radio Communications (PIMRC), 2015 IEEE 26th Annual International Symposium on*, 2015, pp. 1740–1745.
- [44] X. Xia, F. Chen, Q. He, J. Grundy, and H. Jin, "Online collaborative data caching in edge computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 2, pp. 281–294, 2020.
- [45] A. Gharaibeh, A. Khreishah, B. Ji, and M. Ayyash, "A provably efficient online collaborative caching algorithm for multicell-coordinated systems," *IEEE Transactions on Mobile Computing*, vol. 15, no. 08, pp. 1863–1876, 2016.
- [46] R. Kaewpuang, D. Niyato, P. Wang, and E. Hossain, "A framework for cooperative resource management in mobile cloud computing," *IEEE JSAC*, vol. 31, no. 12, pp. 2685–2700, 2013.
- [47] R. Yu, J. Ding, S. Maharjan, S. Gjessing, Y. Zhang, and D. Tsang, "Decentralized and optimal resource cooperation in geo-distributed mobile cloud computing," *IEEE Transactions on Emerging Topics in Computing*, vol. 6, no. 1, pp. 72–84, 2018.



**Deke Guo** received the B.S. degree in industry engineering from the Beijing University of Aeronautics and Astronautics, Beijing, China, in 2001, and the Ph.D. degree in management science and engineering from the National University of Defense Technology, Changsha, China, in 2008. He is currently a Professor with the College of System Engineering, National University of Defense Technology, and a Professor with the School of Computer Science and Technology, Tianjin University. His research interests include distributed systems, software-defined networking, data center networking, wireless and mobile systems, and interconnection networks. He is a senior member of the IEEE and a member of the ACM.



**Xin Li** is currently a software engineer at VMware Inc. He obtained his Ph.D. in Computer Engineering from University of California Santa Cruz in 2018. Before that, he received the B.Eng. degree in Communication Engineering from University of Electronic Science and Technology of China and M.S. degree in Electrical Engineering from University of California Riverside. His research interests include network security, IoT, SDN/NFV, distributed systems.



**Junjie Xie** received the B.E. degree in computer science and technology from the Beijing Institute of Technology, Beijing, China, in 2013. He received the M.E. and Ph.D. degrees in management science and engineering from the National University of Defense Technology, Changsha, China, in 2015 and 2020, respectively. He is currently an engineer with the institute of systems engineering, AMS, PLA, Beijing, China. His research interests include distributed systems, software-defined networking and mobile edge

computing.



**Ge Wang** is an Assistant Professor at Xian Jiaotong University. She received her Ph.D degree at Xian Jiaotong University in 2019. She was a visiting student at University of California, Santa Cruz from 2017 to 2019. Her research interests include wireless sensor network, RFID and mobile computing.



**Chen Qian** is an Associate Professor at the Department of Computer Science and Engineering, UC Santa Cruz. He received the B.Sc. degree from Nanjing University in 2006, the M.Phil. degree from the Hong Kong University of Science and Technology in 2008, and the Ph.D. degree from the University of Texas at Austin in 2013, all in Computer Science. His research interests include computer networking, data-center networks and cloud computing, Internet of Things, and software defined networks. He has published more than 60 research papers in a number of top conferences and journals including *ACM SIGMETRICS*, *IEEE ICNP*, *IEEE ICDCS*, *IEEE INFOCOM*, *IEEE PerCom*, *ACM UBIComp*, *ACM CCS*, *IEEE/ACM Transactions on Networking*, and *IEEE Transactions on Parallel and Distributed Systems*. He is a member of IEEE and ACM.

published more than 60 research papers in a number of top conferences and journals including *ACM SIGMETRICS*, *IEEE ICNP*, *IEEE ICDCS*, *IEEE INFOCOM*, *IEEE PerCom*, *ACM UBIComp*, *ACM CCS*, *IEEE/ACM Transactions on Networking*, and *IEEE Transactions on Parallel and Distributed Systems*. He is a member of IEEE and ACM.



**Honghui Chen** received the MS degree in operational research and the PhD degree in management science and engineering from the National University of Defense Technology, Changsha, China, in 1994 and 2007, respectively. Currently, he is a professor of College of System Engineering, National University of Defense Technology, Changsha, China. His research interests include information system, cloud computing and Information Retrieval.