

On Optimizing Replica Migration in Distributed Cloud Storage Systems

Amina Mseddi*, Mohammad Ali Salahuddin*[‡], Mohamed Faten Zhani[†], Halima Elbiaze*, Roch H. Glitho[‡]

*Université du Québec À Montréal, Montreal, Quebec, Canada

[†]École de Technologie Supérieure, Montreal, Quebec, Canada

[‡]Concordia University, Montreal, Quebec, Canada

mseddi.amina@courrier.uqam.ca, mohammad.salahuddin@ieee.org,

mfzhani@etsmtl.ca, elbiaze.halima@uqam.ca, glitho@ciise.concordia.ca

Abstract—With the wide adoption of large-scale Internet services and big data, the cloud has become the ideal environment to satisfy the ever-growing storage demand, thanks to its seemingly limitless capacity, high availability and faster access time. In this context, data replication has been touted as the ultimate solution to improve data availability and reduce access time. However, replica placement systems usually need to migrate and create a large number of data replicas over time between and within data centers, incurring a large overhead in terms of network load and availability. In this paper, we propose CRANE, an efficient Replica migration scheme for distributed cloud Storage systems. CRANE complements any replica placement algorithm by efficiently managing replica creation in geo-distributed infrastructures by (1) minimizing the time needed to copy the data to the new replica location, (2) avoiding network congestion, and (3) ensuring a minimal availability of the data. Our results show that, compared to swift (the OpenStack project for managing data storage), CRANE is able to minimize up to 30% of the replica creation time and 25% of inter-data center network traffic, while ensuring the minimum required availability of the data.

I. INTRODUCTION

With the wide adoption of large-scale Internet services and big data, the cloud has become the ultimate resort to cater to the ever-growing demand for storage, providing seemingly limitless capacity, high availability and faster access time. Typically, cloud providers build several large-scale data centers in geographically distributed locations. Then, they rely on data replication as an effective technique to provide fault-tolerance, reduce end-user latency and minimize the amount of data exchanged through the network. As a result, replica management has become one of the major challenges for cloud providers.

In recent years, a large body of work has been devoted to several challenges related replica management in distributed cloud storage systems. A large part of the research efforts were mostly dedicated to replica placement problem, considering different goals such as minimizing storage costs, improving fault-tolerance and access delays [1]–[5]. However, replica placement systems may result in a huge number of data replicas created or migrated over time between and within data centers, incurring large amounts of traffic between data centers. This can be the case in different scenarios: for instance, when a new data center is added to the cloud provider’s infrastructure, when a data center is scaled up or down, when recovering from a disaster or simply when achieving performance or

availability goals, requiring the creation and the relocation of a large number of replicas.

Naturally, several impacts may be expected when such large data bulk transfer of replicas is triggered. These impacts can be summarized as follows:

- As copying data consumes resources (e.g., CPU, memory, disk I/O) at both the source and the destination machines, these nodes will experience more contention for the available capacity, which may slow down other tasks running on them.

- Recent research revealed that traffic exchanged between data centers account for up to 45% of the total traffic in the backbone network connecting them [6]. This ever-growing exchange of tremendous amounts of data between data centers may overload the network, especially when using the same paths or links. This can hurt the overall network performance in terms of latency and also packet loss.

- Replica migration processes are usually distributed and asynchronous as is the case for Swift, the OpenStack project for managing data storage [7]. That is, when a replica is to be relocated or created in a new destination machine, every machine in the infrastructure already storing the same replica will try to copy the data to the new destination. There is no coordination or synchronization between the sending nodes. This will not only lead to unneeded redundancy as the same data is copied from different sources at the same time, but will also further exacerbate the congestion in the data center network.

- Replicas are usually placed in geographically distributed locations, so as to increase data availability over time and reduce user-perceived latency. When a replica have to be created/migrated in a new location, it will not be available until all its content is copied from other existing replicas. If this process takes too long, it might hurt the overall data availability, if the number of available replicas is not sufficient to accommodate all user requests. For instance, in order to ensure availability, Swift ensures that at least two replicas of the data are available at any point in time (according to the default configuration [7]).

In order to alleviate all the aforementioned problems, it is critical to make sure that replicas are created as soon as possible in their new locations without inferring network congestion or high creation time. This requires to carefully

select the source replica from which the data will be copied, the paths through which the data will be sent and the order in which replicas are created.

To address these challenges, we start by formulating the replica migration/creation problem as an Integer Linear Program (ILP). We then propose (CRANE¹) an efficient Replica migration scheme for distributed cloud Storage systems. CRANE is a novel scheme that manages replica creation in geo-distributed infrastructures with the goal of (1) minimizing the time needed to copy the data to the new replica location, (2) avoiding network congestion, and (3) ensuring a minimal availability for each replica. CRANE can be used with any existing replica placement algorithm in order to optimize the time to create and copy replicas and to minimize resources needed to reach the new placement of the replicas. In addition, it ensures that at any point in time, data availability is above a predefined minimum value.

The rest of the paper is organized as follows. Section II surveys the related work on replica placement and migration in the cloud. Section III presents an example illustrating how different replica creation and migration strategy can impact performance metrics. We then formulate the replica creation problem in Section IV. Our proposed solution is presented in Section V. Finally, we provide in Section VI some simulation results that assess the performance of CRANE and compare it to Swift and we conclude in Section VII.

II. RELATED WORK

In this section, we survey relevant works on replica management in the cloud. Several efforts have been devoted to put forward effective solutions for replica placement problem. That is to determine the optimal number and placement of data replicas in order to achieve several objectives such that minimizing hosting costs, reducing access delay to the data, and maximizing data availability [1]–[5], [8]. Once the replica placement algorithm is executed, a new placement of replicas is determined in order to achieve the desired objective. Hence, some existing replicas should be torn down and some new ones should be created across the infrastructure.

In this work, we do not focus on the replica placement but rather on reducing the overhead of migrating from an original placement of replicas to the new one, which should take place right after the execution of the replica placement algorithm. In the last few years, very few proposals have looked at this problem but overlooked many important parameters. For instance, Kari et al. [9] proposed a scheme that tries to find an optimal migration schedule for data in order to minimize the total migration time. They take into account the heterogeneity of storage nodes in terms of the number of simultaneous transfers they can handle. However, they have overlooked availability requirements as well as network-related constraints such as bandwidth limits and propagation delays. Other works [10], [11] proposed different approximation algorithms to solve the problem; however they always aim at minimizing migration times without considering the availability of data

during the migration process. Finally, Swift, the OpenStack project for managing data storage [7], implements a data replica placement algorithm along a replica migration one. As a placement strategy, blocs of data (called hereafter as partitions) are replicated and distributed across the distributed infrastructure according to the as-unique-as-possible algorithm [12], which ensure that partition's replicas are physically stored as far as possible from each other in order to ensure high availability. In terms of replica creation and migration, Swift simply do migrates replicas between data centers without considering the network available capacity. However, it ensures a high availability of the data by allowing only one migration per partition each time interval (usually, one hour) so that only one replica can be missing at a particular point of time. Of course, the 1-hour waiting time for triggering migrations will significantly increase the total time needed to reach the new optimal placement of replicas.

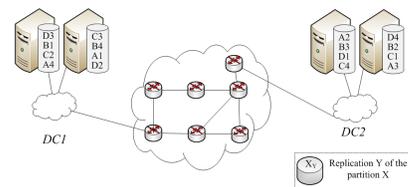
Different from previous work, CRANE takes into account not only the network available capacity but also data availability during the creation of the new replicas. Furthermore, it capitalizes on the existence of multiple replica across the network in order to carefully select the source of the data and the transmission path in order to avoid network congestion and minimize data migration time.

III. MOTIVATING EXAMPLE

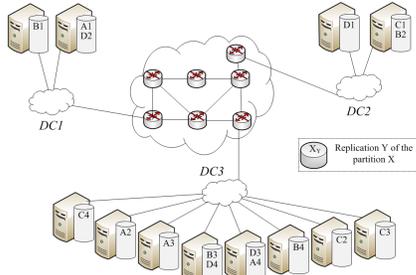
To introduce our proposed replica placement solution, a motivating example is described in this section. Let's consider a cloud system composed of two data centers that span multiple geographic regions. The data centers have different storage capacities and are interconnected by different capacity links. This cloud deployment uses Swift as a distributed solution for managing storage functionalities. Consider a scenario where 4 partitions A, B, C and D with sizes 300 GB, 100 GB, 500 GB and 200 GB, respectively, are configured. Each partition is replicated 4 times and the replicas are stored across data centers. Fig. 1(a) illustrates the initial mapping of the replicas on the data centers.

When a new data center is added, partitions are relocated according to the as-unique-as-possible algorithm, respecting the capacity of the disks on servers in the data centers. To reach the new configuration, bulk volumes of replications will be migrating. Starting from the time where we decide to relocate the replications, all the Swift components will operate according to the new mapping of the replications, directing client requests to non-existent replicas, resulting in unavailability. Swift requires a majority of replicas responding to consider a retrieval operation successful. Thus, to avoid unavailability of the replicas, Swift specifies in its configuration the number of hours to restrict moving a partition more than once. In our scenarios, we assume this parameter to be one hour. Thus, Swift will not move two replicas of the same partition at the same time, but wait an hour before moving the second replica. It is important to note that the second migration is not triggered automatically. Instead, the cloud administrator should run a command to recompute the location of the replications, and then migration starts. Also, in this scenario, the minimum tolerated availability of replicas of each partition is assumed to be $\frac{3}{4}$. Fig. 1(b) shows the optimal locations

¹CRANE:(Mechanical Engineering) a device for lifting and moving heavy objects, typically consisting of a moving boom, beam, or gantry from which lifting gear is suspended



(a) Initial replication mapping



(b) Final replication mapping

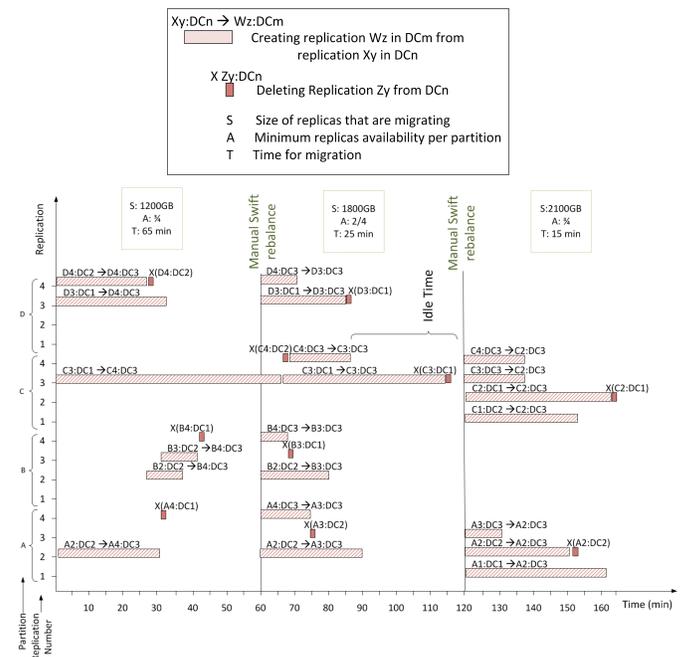
Fig. 1. Initial and final replication mapping

of the different replications according to the disk capacities and the as-unique-as-possible algorithm. The depicted final configuration is reached after several calculations of locations and migrations.

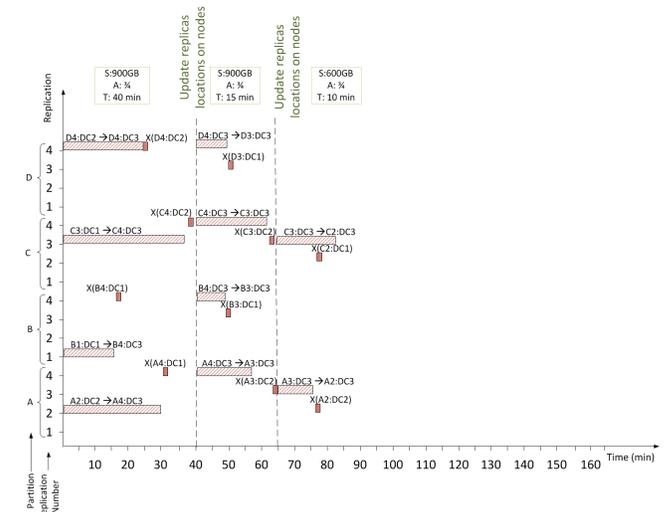
Migration of replications in Swift is a peer-to-peer asynchronous process. Each storage server will sequentially run through its replications and compare them to the ones on the servers assigned to hold a copy of this replication. The migration process uses a push model, so each time a replication needs to be present in one server and it's not, this replication will be copied from the local server to the remote one. And if a server holds a copy of one replication that shouldn't be there it should delete it. Fig. 2(a) represents an example of the replication migration sequence. In this figure, the bottom line represents time in minutes and the replication identifiers of each partition are represented on the vertical line. For instance each storage server starts sequentially copying the replicas to the servers that should hold a copy of them and are not. However, the following problems arise:

Availability: The number of hours to restrict moving a replica of the same partition is a parameter defined by the storage provider. In this scenario it's set to one hour. But migration of replications might take more than one hour. This may have a bad impact on the minimum tolerated availability of replicas per partition. In the depicted scenario, copying replica C3 from DC1 to create C4 on DC3 have taken more than one hour, and the new computation of location of replicas dictated that D3 should be migrated to DC3. So, in the second migration sequence, we have both C4 and C3 that are not available on DC3 during the five first minutes. This violates the minimum required replicas available per partition.

Redundancy: There is a redundancy on the creation of the replica D4 on DC3. In fact, replication D3 and D4 are copied from DC1 and DC2, respectively. This redundancy implies needless bandwidth consumption, which increases replication migration time. Moreover, the copying of B2 from DC2 to



(a) Example of Swift replication migration sequence



(b) Novel replication migration sequence

Fig. 2. Replication migration sequences

DC3 is delayed. Indeed, all the storage servers were copying the first replica in their list, that is, B2 started to be copied only when the server was finished copying D4.

Migration time: Each storage server on Swift starts copying the replicas of the partitions that need to be migrated without any calculation on the delays that this task would take. Therefore, one storage server could end up copying a replica, even though, there exists another server with a replica of the same partition, which could have performed the replication faster. For example, in order to create D4 in DC3, copying D3 from DC1 takes less time than copying D4 from DC2. This could be due to the difference on the available bandwidth on the links and the propagation delays between data centers.

Moreover copying a replica within the same data center will lead to better migration time.

Idle time: We notice there is always an idle time between two sequences of migration.

To avoid the aforementioned drawbacks, we propose an enhanced migration solution illustrated in Fig. 2(b). Note, avoiding the redundancy of copying replicas increases available bandwidth and thus decreases migration time. Also, smart selection of copying source further reduces the migration time significantly. In fact, creating replicas from within the same data center or from a distant data center with links having better propagation delay and less bandwidth consumption will lead to a faster migration time. Furthermore, the minimum replica availability tolerated per partition should be maintained. Finally, automating the recalculation of the location after each sequence of migration leads to avoiding idle time and faster convergence to the final mapping of replicas in data centers.

IV. THE REPLICA MIGRATION PROBLEM

A. Problem Statement

Given a network represented by a graph $G = (S, E)$, where $S = \{s_1, s_2, \dots, s_i, \dots, s_k, \dots, s_{|S|}\}$ is the set of all servers across data centers. We assume that data centers are connected through a backbone network. The backbone's links are represented by a set of edges E . Each edge $e \in E$ is characterized by a bandwidth capacity B_e . Let $P = \{p_1, p_2, \dots, p_j, \dots, p_{|P|}\}$ denote the set of partitions, replicas of which are stored across the servers where $|p_j|$ is the size of replica of partition p_j . We define a configuration as a particular placement of the replicas of partitions within servers. Given an initial and a final configuration, denoted respectively by C^I and C^F , our goal is to find the optimal sequence of replica migrations that minimizes the migration time from C^I to C^F while meeting the minimum partition availability threshold A and abiding by edge bandwidth capacities. We model this as an Integer Linear Programming (ILP) problem. Tables I and II show respectively the inputs of the ILP and its variables.

TABLE I. PROBLEM INPUTS

Input	Definition
S	Set of servers across data centers, where $S = \{s_1, s_2, \dots, s_i, \dots, s_k, \dots, s_{ S }\}$
E	Set of edges connecting servers in S
B_e	Bandwidth capacity $\forall e \in E$
P	Set of partitions, where $P = \{p_1, p_2, \dots, p_j, \dots, p_{ P }\}$ and $ p_j $ is the size of partition p_j
C^I	$ S \times P $ matrix representing an initial configuration, where $c_{i,j}^I = \begin{cases} 1, & \text{if replica of } p_j \text{ is stored on } s_i \\ 0, & \text{otherwise} \end{cases}$
C^F	$ S \times P $ matrix representing a final configuration, where $c_{i,j}^F = \begin{cases} 1, & \text{if replica of } p_j \text{ is stored on } s_i \\ 0, & \text{otherwise} \end{cases}$
T	Worst-case migration time
G	$ S \times S \times E $ matrix representing edges used in a path, where $g_{i,k,e} = \begin{cases} 1, & \text{if edge } e \text{ is used in shortest path between } s_i \text{ and } s_k \\ 0, & \text{otherwise} \end{cases}$
β	A big constant

TABLE II. PROBLEM VARIABLES

Variable	Definition
X	$ S \times P \times S \times T$ matrix representing migration sequence, where $x_{i,j,k,t} = \begin{cases} 1, & \text{if } s_i \text{ is migrating replica of } p_j \text{ to } s_k \text{ at time } t \\ 0, & \text{otherwise} \end{cases}$
Y	$ S \times P \times S $ matrix representing a need for partition migration, where $y_{i,j,k} = \begin{cases} 1, & \text{if } s_i \text{ is the provider of replica of } p_j \text{ to } s_k \\ 0, & \text{otherwise} \end{cases}$
Z	$ S \times P \times T$ matrix representing replica placement, where $z_{i,j,t} = \begin{cases} 1, & \text{if } s_i \text{ has replica of } p_j \text{ at time } t \\ 0, & \text{otherwise} \end{cases}$
L	$ E \times T$ matrix representing load $l_{e,t}$ on edge e at time t , where $l_{e,t} \leq B_e$
R	$ S \times P \times S \times T$ matrix, where $r_{i,j,k,t}$ represents the bandwidth allocated for migrating replica of p_j from s_i to s_k at time t
D	$ S \times P $ matrix representing the replicas that are to be deleted at time T , where $d_{i,j} = \begin{cases} 1, & \text{if replica of } p_j \text{ will be deleted from } s_i \\ 0, & \text{otherwise} \end{cases}$
V	$ S \times S \times T$ matrix, where $v_{i,k,t}$ represents the capacity of the path between s_i to s_k at time t
W	A vector of size T , where $w_t = \begin{cases} 1, & \text{if migration is in progress at time } t \\ 0, & \text{otherwise} \end{cases}$

B. Constraints

Given the initial and final configurations, any discrepancy in the configurations necessitates either the migration or the deletion of the partition replicas. Consider that the migration or deletion of the replica is identified by variables $y_{i,j,k}$ and $d_{k,j}$, respectively. Then, if a server s_k has replica of partition p_j in the initial and final configuration and no action is required, then there should be no migration of replica of partition p_j from any server s_i to s_k , that is, $\sum_{i=1}^{|S|} y_{i,j,k} = 0$, and there should be no deletion of replica of partition p_j from server s_k , that is, $d_{k,j} = 0$, as in (1). However, if the replica of partition p_j is present on server s_k in the initial configuration and not in the final configuration, then the partition should be eventually deleted, hence, $d_{k,j}$ is set to 1, with constraint (2). Most importantly, in constraint (3) we capture the need for a replica migration, if the server s_k does not have the replica of partition p_j in the initial configuration. In this case, there should be some server s_i that delivers the replica of partition p_j to server s_k , therefore $\sum_{i=1}^{|S|} y_{i,j,k} = 1$.

$$c_{k,j}^I + \sum_{i=1}^{|S|} y_{i,j,k} - d_{k,j} = c_{k,j}^F \quad \forall 1 \leq k \leq |S|, 1 \leq j \leq |P| \quad (1)$$

$$d_{k,j} \geq c_{k,j}^I - c_{k,j}^F \quad \forall 1 \leq k \leq |S|, 1 \leq j \leq |P| \quad (2)$$

$$\sum_{i=1}^{|S|} y_{i,j,k} \geq c_{k,j}^F - c_{k,j}^I \quad \forall 1 \leq k \leq |S|, 1 \leq j \leq |P| \quad (3)$$

Before we can initiate the migration, we have to identify the servers s_i that hold the replica of partition p_j at time t , in variable $z_{i,j,t}$ in constraint (4). To begin the migration, only those servers s_i can participate in the replica migration that hold a replica of partition p_j in the initial configuration.

$$c_{i,j}^I \leq \beta \cdot z_{i,j,t} \quad \forall 1 \leq i \leq |S|, 1 \leq j \leq |P|, 1 \leq t \leq T \quad (4)$$

The variable $z_{k,j,t}$ that indicates potential sources of replicas that indicates potential sources of replicas is updated in each time instance t , as servers s_k may begin to hold a copy of the replica of partition p_j , due to migration in earlier time

instances, as in constraints (5) and (6). A server holds a copy of the replica when the sum of all the bandwidth allocated, $r_{i,j,k,t'}$ to the migration of that replica of partition p_j from source s_j to destination s_k , in previous time instances $\forall t' < t$, equals the size of the partition p_j . Then, in following time instances, server s_k could potentially participate in the replica migration.

$$|p_j| - \sum_{k=1, i \neq k}^{|S|} \sum_{t'=1}^t r_{i,j,k,t'} \leq \beta \cdot (1 - z_{k,j,t+1})$$

$$\forall 1 \leq i \leq |S|, 1 \leq j \leq |P|, 1 \leq t \leq T-1 \quad (5)$$

$$|p_j| - \sum_{k=1, i \neq k}^{|S|} \sum_{t'=1}^t r_{i,j,k,t'} \geq 1 - z_{k,j,t+1}$$

$$\forall 1 \leq i \leq |S|, 1 \leq j \leq |P|, 1 \leq t \leq T-1 \quad (6)$$

The bandwidth allocated for migration of replica of partition p_j at time t from source s_i to destination server s_k , in variable $r_{i,j,k,t}$, is dependent on the capacity of the shortest path from s_i to s_k , in variable $v_{i,k,t}$, or the remaining size of the partition replica to be migrated in constraint (7).

$$r_{i,j,k,t} = \operatorname{argmin}\left\{ |p_j| - \sum_{k=1, i \neq k}^{|S|} \sum_{t'=1}^{t-1} r_{i,j,k,t'}, v_{i,k,t} \right\}$$

$$\forall 1 \leq i, k, i \neq k \leq |S|, 1 \leq j \leq |P|, 1 \leq t \leq T \quad (7)$$

The capacity of the shortest path between s_i and s_k is inferred from the load on the edges traversed in the path. The edge with the least capacity, bounds the capacity of the path from above. The available path capacity at time t is in constraint (8). The edges can be traversed by multiple paths, that is, multiple paths between different source destination pairs can share common edges. Therefore, the load on an edge, not exceeding edge capacity, is conjured as the sum of all partition replicas migrating in the network across all source and destination servers at time t , on edge e in the shortest path between s_i and s_k , by constraints (9).

$$v_{i,k,t} = \operatorname{argmin}\left\{ (B_e - l_{e,t}) \cdot g_{i,k,e} \mid 1 \leq e \leq |E| \right\}$$

$$\forall 1 \leq i, k, i \neq k \leq |S|, 1 \leq t \leq T \quad (8)$$

$$l_{e,t} = \sum_{i=1}^{|S|} \sum_{j=1}^{|P|} \sum_{k=1}^{|S|} g_{i,k,e} \cdot r_{i,j,k,t} \quad \forall 1 \leq e \leq |E|, 1 \leq t \leq T \quad (9)$$

Once the model has been initialized with potential providers, need for migration and the bandwidth allocated for the migration of partition replicas, we can initiate migration by associating the replica of partition migration indicator variable $x_{i,j,k,t}$ with need for migration of replica of partition p_j from s_i to s_k , in variable $y_{i,j,k}$, in constraints (10) and (11). Consequentially, from (10) and (11), we ensure only sequential migration of replica of partition, since concurrency is set to 1. Furthermore, in constraint (12), we bind the source server s_i , such that, only those servers that hold complete replica of partition p_j can initiate migration.

$$\sum_{t=1}^T x_{i,j,k,t} \leq \beta \cdot y_{i,j,k} \quad \forall 1 \leq i, k, i \neq k \leq |S|, 1 \leq j \leq |P| \quad (10)$$

$$\sum_{t=1}^T x_{i,j,k,t} \geq y_{i,j,k} \quad \forall 1 \leq i, k, i \neq k \leq |S|, 1 \leq j \leq |P| \quad (11)$$

$$\sum_{k=1}^{|S|} x_{i,j,k,t} \leq z_{i,j,t} \quad \forall 1 \leq i \leq |S|, 1 \leq j \leq |P|, 1 \leq t \leq T \quad (12)$$

To ensure continuous sequential migration of replica of partition p_j , from same source server s_i to same destination server s_k for next time instance $t+1$, we set the indicator of migration is progress, in variable $x_{i,j,k,t+1}$, to 1, until the entire replica of the partition has been migrated. This is depicted in constraints (13) and (14).

$$\beta \cdot x_{i,j,k,t+1} \geq |p_j| - \sum_{t'=1}^t r_{i,j,k,t'}$$

$$\forall 1 \leq i, k, i \neq k \leq |S|, 1 \leq j \leq |P|, 1 \leq t \leq T-1 \quad (13)$$

$$x_{i,j,k,t+1} \leq |p_j| - \sum_{t'=1}^t r_{i,j,k,t'}$$

$$\forall 1 \leq i, k, i \neq k \leq |S|, 1 \leq j \leq |P|, 1 \leq t \leq T-1 \quad (14)$$

During the migration process the servers must maintain the minimum availability threshold for each partition with constraint (15).

$$\sum_{i=1}^{|S|} z_{i,j,t} \geq A \quad \forall 1 \leq j \leq |P|, 1 \leq t \leq T \quad (15)$$

The total migration time is extended to include all migrations in progress in constraint (16) and stopping the migration process in constraint (17).

$$w_t \geq x_{i,j,k,t}$$

$$\forall 1 \leq i, k, i \neq k \leq |S|, 1 \leq j \leq |P|, 1 \leq t \leq T \quad (16)$$

$$w_t \geq w_{t+1} \quad \forall 1 \leq t \leq T-1 \quad (17)$$

C. Objective

$$\operatorname{minimize} \left(\sum_{t=1}^T w_t \right) \quad (18)$$

We minimize the total migration time in (18). As the optimization minimizes migration times, it will select source servers for replica migration that reduce migration time, such that it selects source-destination pairs that have minimum overlapping edges in the shortest path, while ensuring sequential migrations, meeting minimum partition availability threshold and abiding by edge bandwidth capacities.

V. SOLUTION DESIGN

In this section we will describe CRANE, our heuristic solution for the replicas migration problem. Given an initial and a target replicas mapping in data centers, the goal of this algorithm is to find the best sources for copying the replicas and the best sequence to send them so as to minimize the total replica creation/migration time. To this end, the following sights can guide the replication replica creation/migration sequence : (1) avoid redundancy, (2) select the source of data and paths having more available bandwidth, and (3) avoid idle time between sequences.

Our heuristic solution is described in Algorithm 1. Given an initial and target placement configurations (i.e., C^I and C^F), Algorithm 1 returns a set Q of sequences Q_i for migration.

Each sequence contains an ordered set of replicas to be migrated/created such that the required minimum availability per partition is satisfied. After each sequence of migrations Q_i , cloud storage components will be updated with the new placement, so that data user requests can be redirected to the right partition locations. The final replica placement is reached once all the sequences $Q_i, i < n$ are executed.

Initially, P contains the set of partitions that need to be created/migrated. This set can be computed based on the initial and the final partition locations (i.e., C^I and C^F). We then initialize Q that should contain the sequence of replicas to be migrated. In line 3, we initialize a variable i that denotes the number of the sequence. The core of the algorithm aims to iteratively add a partition replica on ordered sequence Q_i . We create a new sequence whenever it is not possible anymore to add a replica creation/migration to the current sequence (not possible because otherwise we do not satisfy the minimum replica availability per partition).

The variable *available* is *true* if there still some replica that can be added to the sequence Q_i . As long as *available* is *true*, we iterate over all partitions in P and we choose the replica r_b from the set of replicas R_p of each partition p that minimizes the migration time. For that we use the variable $T_{Q_i, R_p, min}$ that denotes the minimal migration time of the sequence Q_i when a replica of partition p is included. We compare this variable to $T_{Q_i, r}$, the migration time of the sequence Q_i if replica r is included (line 15). This allows us to select the best replica r_b of the partition p . From the selected replicas, we need to select the one that minimizes migration time (denoted by r_s). To do so, we use the variable $T_{Q_i, min}$ which denotes the minimal migration time after adding a new replica to the sequence Q_i (line 9). We then choose from all previously selected replicas the one that minimizes migration time after adding a new replica (line 20 to 23). The chosen replica is then added to Q_i (line 26). The partition p whose replica r_s was selected is then removed from the set P .

To detect that we cannot add any more replica to the sequence Q_i , we iterate over all partition until the variable *available* becomes *True*. At that time, we add the sequence to Q , and start a new one as long as we still have partitions in P to migrate/create.

VI. PERFORMANCE EVALUATION

In this section, we compare the performance of Swift using CRANE for replica migration with the traditional Swift, with respect to migration time, amount of transferred data and partition availability.

A. Deployment scenarios

Our evaluation environment consists of five data centers, each consisting of five storage servers. We use the NSFNet topology as a backbone network interconnecting the data centers [13]. We use the standard Swift configuration stipulating that for each data partition, three replicas have to be created and placed according as-unique-as possible algorithm. Furthermore, swifts assumes that if 2 out of the 3 replicas are available, the data is assumed to be available (i.e., all user

Algorithm 1 CRANE

require: Initial configuration C^I .
require: Final configuration C^F
output: Sequence for migration

```

1:  $P \leftarrow$  partitions to be migrated
2:  $Q \leftarrow \{\emptyset\}$ 
3:  $i \leftarrow 0$ 
4: while  $P \neq \{\emptyset\}$  do
5:    $Q_i \leftarrow \{\emptyset\}$ 
6:   available  $\leftarrow$  True
7:   while available == True do
8:     available  $\leftarrow$  False
9:      $T_{Q_i, min} \leftarrow \infty$ 
10:    for each  $p \in P$  do
11:      if  $p$  can be included in  $Q_i$  then
12:        available  $\leftarrow$  True
13:         $T_{Q_i, R_p, min} \leftarrow \infty$ 
14:        for each  $r \in R_p$  do
15:          if  $T_{Q_i, r} < T_{Q_i, R_p, min}$  then
16:             $T_{Q_i, R_p, min} = T_{Q_i, r}$ 
17:             $r_b = r$ 
18:          end if
19:        end for
20:        if  $T_{Q_i, R_p, min} < T_{Q_i, min}$  then
21:           $T_{Q_i, min} = T_{Q_i, R_p, min}$ 
22:           $r_s = r_b$ 
23:        end if
24:      end if
25:    end for
26:     $Q_i = Q_i \cup r_s$ 
27:     $P = P - \{\text{partition } p \text{ of replica } r_s\}$ 
28:  end while
29:   $Q = Q \cup Q_i$ 
30:   $i \leftarrow i + 1$ 
31: end while
32: return  $Q$ 

```

requests can be accommodated). Hence, the minimum required availability per partition is set to 2/3.

In our simulations, we consider four scenarios as depicted in Table III. Each having different number of partitions placed across the data center with different partition sizes. In the beginning of each experiment, we consider only 4 data centers. After that, a new data center is connected to the infrastructure, which triggers the placement algorithm in order to re-optimize the location of replicas. We then use whether swift combined with CRANE or the traditional swift to migrate or create new replicas. For instance, in scenario 1, we originally have 512 partitions (i.e., 1536 replicas) distributed across the infrastructure. When the fifth data center is added, 656 replicas should be migrated or created across the fifth data centers (Table III).

B. Results

For each scenario, depicted in Table III, we compare (Swift + CRANE) with traditional Swift with respect to the

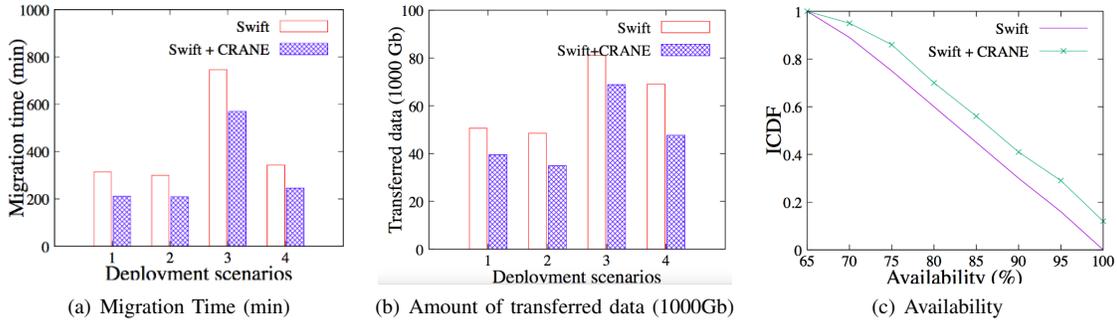


Fig. 3. Performance comparison between CRANE and traditional Swift.

TABLE III. DEPLOYMENT SCENARIOS

Scenario	Total number of partitions	Number of replicas to migrate	Range of replicas size (Gb)
1	512	656	50-100
2	1024	1316	20-50
3	2048	2632	20-50
4	4094	5264	10-20

following performance metrics: (1) the total migration time, (2) the amount of inter-data centers exchanged data, and (3) the availability of replicas per partition.

Figure 3(a) shows the total migration time for the considered scenarios. As we can see, in all scenarios, (CRANE + Swift) outperforms the Swift algorithm by a good margin. For scenario 1, the swift algorithm takes 315 min to create all the replicas compared to 217 min for CRANE, which constitutes around 30% of improvement. For scenario 2, the replicas are migrated within 300 min with swift and 200 min with CRANE, which constitutes around 30% improvement. For scenarios 3 and 4, CRANE achieves 25% improvement. These results are as expected, because CRANE always chooses to copy the replica incurring the minimal transmission time.

The amount of transferred data inter-data centers is reported in figure 3(b). For the 4 different scenarios the CRANE algorithm have less amount of data transferred. The improvement is around 25%. This improvement is explained by the avoidance of redundant copy of the replicas of the same partition. This have also induced the improvement in migration time showed in figure 3(a).

Finally, figure 3(c) shows the Inverse Cumulative Distribution Function (ICDF) of the availability. For a given availability, it provides the probability of having that availability or higher. The required minimum availability per partition ($2/3 = 0.66$) is always met for both algorithms as we can see that the probability of having an availability higher than 66% is 1. However, the probability of having a high availability is always higher for the CRANE algorithm than the traditional Swift. For instance, the probability of having an availability higher than 80% is 0.60 for Swift whereas it is around 0.76 for CRANE. When comparing the two curves, we can see that, on average, CRANE improves availability by up to 10%.

It is clear that CRANE performs significantly better than the basic Swift algorithm as it carefully selects the replica from which the data should be copied, the paths used to transmit that

data while avoiding redundant copy of replicas and eliminating idle time.

VII. CONCLUSION

Data replication has been widely adopted to improve data availability and to reduce access time. However, replica placement systems usually need to migrate and create a large number of replicas between and within data centers, incurring a large overhead in terms of network load and availability. In this paper, we proposed CRANE, an efficient Replica migration scheme for distributed cloud Storage systems. CRANE complements replica placement algorithms by efficiently managing replica creation by minimizing the time needed to copy data to the new replica location while avoiding network congestion and ensuring the required availability of the data. In order to evaluate the performance of CRANE, we compare it to the standard swift, the OpenStack project for managing data storage. Simulations show that CRANE is able to reduce up to 30% of the replica creation time and 25% of inter-data center network traffic and provide better data availability during the process of replica migration.

REFERENCES

- [1] S. Ghemawat, H. Gobioff, and S.-T. Leung, "The Google file system," in *ACM SIGOPS operating systems review*, vol. 37, no. 5, 2003.
- [2] K. Shvachko, H. Kuang, S. Radia, and R. Chansler, "The hadoop distributed file system," in *IEEE Symposium on Mass Storage Systems and Technologies (MSST)*, 2010.
- [3] R.-S. Chang and H.-P. Chang, "A dynamic data replication strategy using access-weights in data grids," *The Journal of Supercomputing*, vol. 45, no. 3, 2008.
- [4] Q. Wei, B. Veeravalli, B. Gong, L. Zeng, and D. Feng, "CDRM: A cost-effective dynamic replication management scheme for cloud storage cluster," in *IEEE International Conference on Cluster Computing (CLUSTER)*, 2010, pp. 188–196.
- [5] D.-W. Sun, G.-R. Chang, S. Gao, L.-Z. Jin, and X.-W. Wang, "Modeling a dynamic data replication strategy to increase system availability in cloud computing environments," *Journal of Computer Science and Technology*, vol. 27, no. 2, pp. 256–272, 2012.
- [6] Y. Chen, S. Jain, V. Adhikari, Z.-L. Zhang, and K. Xu, "A first look at inter-data center traffic characteristics via Yahoo! datasets," in *IEEE INFOCOM*, April 2011, pp. 1620–1628.
- [7] O. foundation. (2015) Swift documentation. [Online]. Available: <http://docs.openstack.org/developer/swift/>
- [8] S. Zaman and D. Grosu, "A distributed algorithm for the replica placement problem," *IEEE Transactions on Parallel and Distributed Systems*, vol. 22, no. 9, Sept 2011.

- [9] C. Kari, Y.-A. Kim, and A. Russell, "Data migration in heterogeneous storage systems," in *IEEE International Conference on Distributed Computing Systems (ICDCS)*, 2011, pp. 143–150.
- [10] S. Khuller, Y.-A. Kim, and Y.-C. J. Wan, "Algorithms for data migration with cloning," in *ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, ser. PODS, 2003.
- [11] E. Anderson, J. Hall, J. Hartline, M. Hobbs, A. Karlin, J. Saia, R. Swaminathan, and J. Wilkes, "An experimental study of data migration algorithms," in *Algorithm Engineering*, ser. Lecture Notes in Computer Science, 2001, vol. 2141.
- [12] J. Dickinson. (2013) Data placement in swift. [Online]. Available: <https://swiftstack.com/blog/2013/02/25/data-placement-in-swift/>
- [13] NSFNET. (2015) National science foundation network (nsfnet). [Online]. Available: <http://hpwren.ucsd.edu/~hwb/NSFNET/NSFNET-200711Summary/>