

Dear Author/Editor,

Here are the proofs of your chapter as well as the metadata sheets.

# Metadata

- Please carefully proof read the metadata, above all the names and address.
- In case there were no abstracts for this book submitted with the manuscript, the first 10-15 lines of the first paragraph were taken. In case you want to replace these default abstracts, please submit new abstracts with your proof corrections.

# Page proofs

- Please check the proofs and mark your corrections either by
  - entering your corrections online or
  - opening the PDF file in Adobe Acrobat and inserting your corrections using the tool "Comment and Markup"
    - or
  - printing the file and marking corrections on hardcopy. Please mark all corrections in dark pen in the text and in the margin at least ¼" (6 mm) from the edge.
- You can upload your annotated PDF file or your corrected printout on our Proofing Website. In case you are not able to scan the printout , send us the corrected pages via fax.
- Please note that any changes at this stage are limited to typographical errors and serious errors of fact.
- If the figures were converted to black and white, please check that the quality of such figures is sufficient and that all references to color in any text discussing the figures is changed accordingly. If the quality of some figures is judged to be insufficient, please send an improved grayscale figure.

# Metadata of the chapter that will be visualized online

Book Title		Handbook on Data Centers
Chapter Title		Cloud Storage over Multiple Data Centers
Copyright		Springer Science+Business Media New York 2014
Corresponding Author [Aff 1]	Family name	Ми
	Particle	
	Given name	Shuai
	Suffix	
	Division	Department of Computer Science and Technology
	Organization	Tsinghua National Laboratory for Information Science and Technology (TNLIST), Tsinghua University
	Address	Beijing, China
	email	msmummy@gmail.com
Author [Aff 2]	Family name	Ми
	Particle	
	Given name	Shuai
	Suffix	
	Division	
	Organization	Research Institute of Tsinghua University in Shenzhen
	Address	Shenzhen, China
	email	msmummy@gmail.com
Author	Family name	Su
	Particle	
	Given name	Maomeng
	Suffix	
	Division	
	Organization	Tsinghua University
	Address	Beijing, China
	email	maomengsu19881010@gmail.com
Author	Family name	Gao
	Particle	
	Given name	Pin
	Suffix	
	Division	
	Organization	Tsinghua University
	Address	Beijing, China
	email	pin.gao2008@gmail.com
Author	Family name	Wu
	Particle	
	Given name	Yongwei
	Suffix	
	Division	
	Organization	Tsinghua University
	Address	Beijing, China
	email	wuyw@tsinghua.edu.cn
Author	Family name	Li
	Particle	
	Given name	Keqin
	Suffix	
	Division	Department of Computer Science

	Organization	State University of New York at New Paltz
	Address	New Paltz, USA
	email	lik@newpaltz.edu
Author	Family name	Zomaya
	Particle	
	Given name	Albert Y.
	Suffix	
	Division	Centre for Distributed and High Performance Computing School of Information Technologies
	Organization	The University of Sydney
	Address	Sydney, Australia
	email	albert.zomaya@sydney.edu.au
Abstract		Cloud storage has become a booming trend in the last few years. Individual developers, companies, organizations, and even governments have either taken steps or at least shown great interests in data migration from self-maintained infrastructure into cloud.

Shuai Mu, Maomeng Su, Pin Gao, Yongwei Wu, Keqin Li and Albert Y. Zomaya

### 1 **Introduction**

- Cloud storage has become a booming trend in the last few years. Individual devel [AQ1]
   opers, companies, organizations, and even governments have either taken steps or at
   least shown great interests in data migration from self-maintained infrastructure into
   cloud.
- 6 Cloud storage benefits consumers in many ways. A recent survey among over 600
- 7 cloud consumers [80] has shown that primary reasons for most clients in turning to
- <sup>8</sup> cloud are: (1) to have highly reliable as well as available data storage services; (2) to
- <sup>9</sup> reduce the capital cost of constructing their own datacenter and then maintaining it;

S. Mu (🖂)

Department of Computer Science and Technology, Tsinghua National Laboratory for Information Science and Technology (TNLIST), Tsinghua University, Beijing, China Research Institute of Tsinghua University in Shenzhen, Shenzhen, China e-mail: msmummy@gmail.com

M. Su · P. Gao · Y. Wu Tsinghua University, Beijing, China e-mail: maomengsu19881010@gmail.com

P. Gao e-mail: pin.gao2008@gmail.com

Y. Wu e-mail: wuyw@tsinghua.edu.cn

K. Li Department of Computer Science, State University of New York at New Paltz, New Paltz, USA e-mail: lik@newpaltz.edu

A. Y. Zomaya Centre for Distributed and High Performance Computing School of Information Technologies, The University of Sydney, Sydney, Australia e-mail: albert.zomaya@sydney.edu.au

© Springer Science+Business Media New York 2014 S. U. Khan, A. Y. Zomaya (eds.), Handbook on Data Centers, DOI 10.1007/978-1-4939-2092-1\_24 1

and (3) to provide high-quality and stable network connectivity to their customers.
 The prosperity of the cloud market has also inspired many companies to provide
 cloud storage services of different quality to vast variety of companies.

13 Reliability and availability are the most important issues in designing a cloud storage system. In cloud storage, data reliability often refers to that data is not 14 lost, and availability refers to data accessibility. To major cloud storage providers, 15 accidents of data loss seldom happen. But there was an accident that Microsoft 16 once completely lost user data of T-Mobile cellphone users [81]. Because data loss 17 accidents seldom happen, data availability is often a more important concern to 18 most cloud storage consumers. Almost all cloud storage providers have suffered 19 from temporary failures, lasting from hours to days. For example, Amazon failed to 20 provide service in October, 2012, which caused many consumers such as Instgram 21 to halt their service. 22

Data replication is an effective way to improve reliability and availability. Limited by cost, cloud providers usually use commodity hardware to store consumers' data. Replicating data into different machines can tolerate hardware failures, and replicating data into multiple data centers can tolerate failures of a datacenter, caused by earthquakes, storms and even wars.

To reduce cost of replication, data are usually divided into stripes instead of the original copy. There are many coding methods for this, originating from traditional storage research and practice. Erasure coding, which is widely used in the RAID systems, provides suitable features for data striping requirements in cloud storage environment.

Data consistency is also an important issue in cloud storage. As opposed to traditional storage systems which usually provide a strong consistency model, cloud storage often offers weaker consistency model such as eventual consistency. Some also propose other reduced consistency models such as session consistency, and fork-join consistency.

Privacy and security are very essential to some consumers. Consumers are often concerned about how their data are visible to cloud providers; whether administrators can see their data transparently. For other consumers who are less sensitive to privacy, they are more concerned about access control to data, because all the traffic finally leads to charges.

In spite of all the efforts of cloud storage providers, there is an emerging trend 43 to build an integration layer on top of current cloud storages, also named "cloud-of-44 clouds". A cloud-of-cloud system makes use of current cloud storage infrastructure, 45 but still provides a uniformed user interface to top-level application developers. It also 46 targets the reliability, availability and security issues, but takes a different approach 47 of using each cloud as a building block. The advantage of this approach is that it 48 can tolerate performance fluctuation of single cloud storage, and can avoid potential 49 risks of a provider's shutdown. 50

In the remainder of this chapter, we first review cloud storage architecture at a high level in Sect. 2. Section 3 describes common strategies used in data replication.

Section 4 gives a brief introduction on data striping. Section 5 introduces the consis tency issue. Section 6 briefly highlights a new model of cloud-of-clouds. Section 7
 discusses privacy and security issues in storage systems. Section 8 summarizes and
 suggests future directions to cloud storage research and practice.

### 57 2 Cloud Storage in a Nutshell

In this section we give an overview of cloud storage architecture and its key components. Cloud storage environments are usually complex systems mixed with many
mature and new techniques. On a high level, cloud storage design and implementation
consist of two parts: metadata and data, which will be investigated later.

### 62 2.1 Architecture

Early public clouds and most of today's private ones are built into a single datacenter,
or several datacenters in nearby buildings. They are composed of hundreds or even
thousands of commodity machines and storage devices, connected by high-speed
networks. Besides large amounts of hardware, many other storage middleware such
as distributed file systems are also necessary to provide storage service to consumers.
The typical architecture of cloud storage usually includes storage devices, distributed
file system, metadata service, frontend, and other components.

In practice, we find that the data models and library interfaces of different clouds
are fairly similar; thus, we could support a minimal set to satisfy most users' needs.
The data model shared by most services could be summarized as a "container-object"
model, in which file objects are put into containers. Most services containers do not
support nesting; i.e., users cannot create a sub-container in a container.

In the last decade, major cloud storage such as Amazon S3 [1] and Windows Azure Storage [2] have upgrade their service from running in separate datacenters to different data centers and different geographic regions. Compared to a single datacenter structure, running services in multiple data centers require a multitude of more resource management functions, such as, resource allocation, deployment, and migration.

An important feature of cloud storage is the ability to store and provide access to an immense amount of storage. Amazon S3 currently has a few hundred petabytes of raw storage in production, and it also has a few hundred more petabytes of raw storage based on customer demand [21]. Modern cloud storage architecture could be divided into three layers: storage service, metadata service, and front-end layer (Fig. 1).

Metadata service—The metadata service is in charge of following functions: (a)
 handling high level interfaces and data structures; (b) managing a scalable names pace for consumers' objects; (c) storing object data into the storage service.



Fig. 1 Cloud storage architecutre over mutiple data centers

Metadata service holds the responsibility to achieve scalability by partitioning 90 all of the data objects within a datacenter. This layer consists of many meta-91 data servers, each of which serves for a range of different objects. Also, it should 92 provide load balance among all the metadata servers to meet the traffic of requests. 93 Storage service—This storage service is in charge of storing the actual data into 94 disks and distributing and replicating the data across many servers to keep data 95 reliable within a datacenter. The storage service can be thought of as a distributed 96 file system. It holds files, which are stored as large storage chunks. It also un-97 derstands how to store them, how to replicate them, and so on, but it does not 98 understand higher level semantics of objects. The data is stored in the storage 99 service, but it is accessed from the metadata service. 100

Front-End (FE) layer—The front-end layer consists of a set of stateless servers that 101 take incoming requests. Upon receiving a request, an FE looks up the account, 102 authenticates and authorizes the request, then routes the request to a partition 103 server in the metadata service. The system maintains a map that keeps track of 104 the partition ranges and which metadata server is serving which partition. The FE 105 servers cache the map and use it to determine which metadata server to forward 106 each request to. The FE servers also file large objects directly from the storage 107 service and cache frequently accessed data for efficiency. 108



Fig. 2 Meta-layer architecture

### 109 2.2 Metadata Service

110 The metadata service contains three main architectural components: a layout man-

ager, many meta-servers, and a reliable lock service (Fig. 2). The architecture is similar to Bigtable [5].

### 113 2.2.1 Layout Manager

A layout manager (LM) acts as a leader of the meta-service. It is responsible for 114 dividing the whole metadata into ranges and assigning each meta-server to serve 115 several ranges and then keeping track of the information. The LM stores this assign-116 ment in a local map. The LM must ensure that each range is assigned only to one 117 active meta-server, and that two ranges do not overlap. It is also in charge of load 118 balancing ranges among meta-servers. Each datacenter may have multiple instances 119 of the LM running, but usually they function as reliable replications of each other. 120 For this they need a Lock Service to maintain a lease for leader election. 121

### 122 2.2.2 Meta-Server

A meta-server (MS) is responsible for organizing and storing a certain set of ranges of metadata, which is assigned by LM. It also serves requests to those ranges. The MS stores all metadata into files persistently on disks and maintains a memory cache for efficiency. Meta-servers keep leases with the Lock Service, so that it is guaranteed that no two meta-servers can serve the same range at the same time.

<sup>128</sup> If a MS fails, LM will assign a new MS to serve all ranges served by the failed

<sup>129</sup> MS. Based on the load, LM may choose a few MS rather than one to serve the ranges.

LM firstly assigns a range to a MS, and then updates its local map which specifies



Fig. 3 Storage service architecture

which MS is serving each range. When a MS gets a new assignment from LM, it firstly acquires for the lease from Lock Service, and then starts serving the new range.

### 133 2.2.3 Lock Service

Lock Service (LS) is used by both of layout manager and meta-server. LS uses Paxos
[16] protocol to do synchronous replication among several nodes to provide a reliable
lock service. LM use LS for leader election; MS also maintains a lease with the LS
to keep alive. Details of the LM leader election and the MS lease management are
discussed here. We also do not go into the details of Paxos protocol. The architecture
of lock service is similar to Chubby [18].

### 140 2.3 Storage Service

The two main architecture components of the storage service are the namenode and chunk server (Fig. 3). The storage service architecture is similar to GFS [4].

#### 143 **2.3.1 Namenode**

The namenode can be considered as the leader of the storage service. It maintains file namespace, relationships between chunks and each file, and the chunk locations across the chunk servers. The namenode is off the critical path of client read and write requests. In addition, the namenode also monitors the health of the chunk servers periodically. Other functions of namenode include: lazy re-replication of chunks, garbage collection, and erasure code scheduling.

The namenode periodically checks the state of each chunk server. If the namenode finds that the replication number of a chunk is smaller than configuration, it will start a re-replication of the chunk. To achieve a balanced chunk replica placement, the namenode randomly chooses chunk server to store new chunk.

6

The namenode is not tracking any information about blocks. It remembers just files and chunks. The reason of this is that the total number of blocks is so huge that the namenode cannot efficiently store and index all of them. The only client of data service is the metadata service.

#### 158 2.3.2 Chunk Servers

Each chunk server keeps the storage for many chunk replicas, which are assigned 159 by the namenode. A chunk server machine has many large volume disks attached, 160 to which it has complete access. A chunk server deals only with chunks and blocks, 161 and it does not care about file namespace in the namenode. Internally on a chunk 162 server, every chunk on disk is a file consisting of data blocks and their checksum. A 163 chuck server also holds a map which specifies relationships between chunk and file. 164 Each chunk server also keeps a view about the chunks it owns and the location of 165 the peer replicas for a given chunk. This view is copied from namenode and is kept 166 as a local cache by the chunk server. Under instructions from namenode, different 167 chunk servers may talk to each other to replicate chunks, or to create new copies of 168 an existing replica. When a chunk no longer stores any alive chunks, the namenode 169 starts garbage collection to remove the dead chunks and free the space. 170

### **3 Replication Strategies**

### 172 3.1 Introduction

Currently, more data-intensive applications are moving their large-scale datasets 173 into cloud. To provide high availability and durability of storage services as well 174 as improving performance and scalability of the whole system, data replication is 175 adopted by many mature platforms [1, 2, 4, 6, 12] and research studies [7-10, 14, 30]176 in cloud computing and storage. Data replication is to keep several identical copies of 177 a data object in different servers that may distribute across multiple racks, houses and 178 region-scale or global-scale datacenters, which can tolerate different levels of failures 179 such as facility outages or regional disasters [4, 10, 23, 30]. Replication strategy is 180 now an indispensable feature in multiple datacenters [1, 2, 6-9, 12], which may 181 be hundreds or thousands of miles away from each other, to completely replicate 182 data objects of services, not only because wide-area disasters such as power outages 183 or earthquakes may occur in one datacenter [10, 23], but also because replication 184 across geographically distributed datacenters can mostly reduce latency and improve 185 the whole throughput of the services in the cloud [6-9, 11]. 186

Availability and durability is guaranteed as one data object is replicated on many servers across datacenters, thus in the presence of failing of a few number of components such as servers and network at any time [1, 4, 10, 23] or natural disasters occurring in one datacenter, the durable service of cloud storage won't be influenced

because applications can normally access their data through servers containing replicas in other datacenters. Moreover, as each data object is replicated over multiple
datacenters, it enables different applications to be served from the fastest datacenter
or the datacenter with the lowest working load in parallel [1, 6, 9, 11, 31], thus
providing high performance and throughput of the overall cloud storage system.

Common replication strategies can be divided into two categories: asynchronous 196 replication and synchronous replication. They own distinct features and have different 197 impacts on availability, performance, and throughput of the whole system. Besides, 198 the cloud storage service should provide the upper applications with a consistent view 199 of the data replicas especially during faults [6-9, 11], which requires that data copies 200 among diverse datacenters should be consistent with each other. However, these two 201 replication strategies bring in new challenges to replication synchronization, which 202 finally will influence the consistency of data replicas over multiple datacenters. 203

Additionally, the placement of data replicas is also an important aspect of replication strategy in multiple datacenters as it highly determines the load distribution, storage capacity usage, energy consumption and access latency, and many current systems and studies [1, 4, 6, 10, 24, 26] adopt different policies for the placement of data replicas in the multiple-datacenter design on different demands.

In this section, we will present the main aspects and features of asynchronous replication, synchronous replication, and placement of replicas.

### 211 3.2 Asynchronous Replication

Figure 4a illustrates the working mechanism of asynchronous replication over mul-212 tiple datacenters. As shown in Fig. 4, the external applications issue write requests 213 to one datacenter, which could be a fixed one configured previously or a random one 214 chosen by applications, and get a successful response if the write requests completely 215 commit in this datacenter. The updated data will be eventually propagated to other 216 datacenters in background in an asynchronous manner [1, 2, 12]. Asynchronous 217 replication is especially useful when the network latency between datacenters is at a 218 high cost as applications only need to commit their write requests in one fast datacen-219 ter and don't have to wait for the data to be replicated in each datacenter. Therefore, 220 the overall performance and throughput for writes will be improved and systems with 221 asynchronous replication can provide high scalability as they are decentralized. Now 222 many prevailing systems such as Cassandra [12], Dynamo [1], and PNUTS [14] are 223 using asynchronous replication. 224

However, asynchronous replication presents a big challenge to consistency, since replicas may have conflicting changes with each other, that is, the view of all the replicas over multiple datacenters has the probability to be inconsistent at some time. Figure 5 presents a simple scenario that will cause inconsistency among replicas. Assume there are three datacenters A, B and C, and all of them hold data replica d. When a write request for d from application P is issued to A and successfully commits in A, A will response to P and then replicates the updated data  $d_1$  to B and C. However,

 $D_l$  for one datacenter



Fig. 4 The working mechanism of asynchronous replication and synchronous replication over multiple datacenters. **a** for asynchronous replication and **b** for synchronous replication



Fig. 5 A scenario that causes inconsistent view of data replicas among datacenters under asynchronous replication

at the same time, another write request for *d* from application *Q* is issued to *C*. As *C* hasn't gotten to know the update of *d* in *A*, it normally accepts and processes this write request and then *d* in *C* turns into  $d_2$  and will be replicated to *A* and *B*. As a result, there are now two different versions of the same data replica, and the system steps into an insistent state which means that a subsequent read may get two different data objects. As there are also other factors such as server or network failure that will cause inconsistency in asynchronous replication over multiple datacenters [1, 11], a few

researches have been addressing this challenge of asynchronous replication. Eventual 239 consistency [21] model is one scheme that is widely adopted by many studies and 240 widespread distributed systems [1, 12-14]. Eventual consistency model allows the 241 whole system to be inconsistent temporarily but eventually, the conflicted data objects 242 will merge into one singe data object and the view of the data replicas across multiple 243 datacenters will become consistent at last. The process of merging conflicted data 244 objects is critical in eventual consistency model and the merging decision can be made 245 by the write timestamp [12, 21], a chosen master [13, 14] or even the applications [1]. 246 Under asynchronous replication, a read request may get a stale data object from 247 some datacenters, which will decline the performance of current reads and com-248 plicates application development. However, whether this circumstance is adverse 249 depends on the applications. If applications such as search engine and shopping 250

carts allow weaker consistency at reading or demand high quality of writing experience [1, 12], asynchronous replication won't bring negative impacts to these applications.

### 254 3.3 Synchronous Replication

In contrast to asynchronous replication, synchronous replication requires that the 255 updated data objects of write requests must be synchronously replicated to all or a 256 majority of datacenters before applications get a successful response from the data-257 center accepting the requests in the cloud, as presented in Fig. 4b. This synchronous 258 replication mechanism can effectively guarantee a consistent view of cross-datacenter 259 replicated data and it enables developers to build distributed systems that can provide 260 strong consistency and a set of ACID semantics like transactions, which, com-261 pared with that in loosely consistent asynchronous replication, simplifies application 262 building for the wide-area usage for the reason that applications can make use of se-263 rializable semantic properties while are free from write conflicts and system crashes 264 [6, 9, 11, 20, 25]. 265

The key point of synchronous replication is to keep states of replicas across 266 different datacenters the same. A simple and intuitive way to realize this is to use 267 synchronous master/slave mechanism [4, 6, 11]. The master waits for the writes to 268 be fully committed in slaves before acknowledging to applications and is responsible 269 for failure detection of the system. Another method to maintain consistent and up-270 to-date replicas among datacenters is to use Paxos [16], which is a fault-tolerant and 271 optimal consensus algorithm for RSM [15] in a decentralized way. Paxos works well 272 when a majority of datacenters are alive and at current, many system services adopt 273 Paxos [2, 6, 11, 17, 18, 20] as their underlying synchronous replication algorithm. 274 However, no matter which method is used, the overall throughput and perfor-275

mance of the services based on synchronous replication will be constrained when the communication latencies between datacenters are at high expense [7, 9, 11] and scalability is limited by strong consistency to certain extent. As a result, many researches put forward mechanisms to help improve the throughput and scalability of

the whole system while not destroying the availability and consistency for applications. These mechanisms include reasonable partitioning of data [2, 6, 11], efficient
use of concurrent control [7, 11, 31] and adopting combined consistent models [7–9,
22, 25].

### 284 3.4 Placement of Replicas

As cloud storage now holds enormous amount (usually petabytes) of data sets from large-scale applications, how to place data replicas across multiple datacenters also becomes a very important aspect in replication strategy as it is closely related to load balance, storage capacity usage, energy consumption and access latency of the whole system [6, 10, 19, 24]. It is essential for both efficiently utilizing available datacenter resources and maximizing performance of the system.

Unbalanced replica placement will cause over-provisioning capacity and skewed 291 utilization of some datacenters [26]. One way to address this issue is to choose 292 a master or use partition layers to decide in which datacenter each data replica is 293 placed [2, 4, 6]. This requires the master or partition layers to record the up-to-date 294 load information of every datacenter so that they won't make unbalanced replica 295 placement policies and can immediately decide to migrate data between datacenters 296 to balance load. Another way is to use a decentralized method, as presented in Fig. 6. 297 We can form datacenters as a ring, each responsible for a range of keys. A data 298 object can get its key through hash functions such as consistent hash and locate a 299 datacenter according to its key. Then, replicas of this data object could be placed in 300 this datacenter and its successive ones, similar to [1, 12]. In this way, there is no need 301 to maintain a master to keep information of each datacenter and if the hash functions 302 could evenly distribute the keys, load balance can be achieved automatically. 303

Furthermore, as datacenters now consumes about 1.5 % of the world's total energy 304 and a big fraction of it does come from the consumption of storage in them [28,305 29], the number of datacenters to place the data replicas should also be considered 306 carefully. If the number of datacenters to hold replicas increases, the storage capacity 307 of the whole system will accordingly decease and the energy consumption improves 308 [24, 26, 27] as those datacenters will contain large amounts of replicated data objects 309 in storage. In addition, placing data replicas in a high number of datacenters enables 310 applications to survive wide-area disasters that will cause a few datacenter failures 311 and thus, this can provide high availability for applications at the expense of storage 312 capacity and energy consumption [2, 6, 7, 11, 25]. Moreover, when the number 313 of datacenters to place replicas is large, applications can have a low access latency 314 based on geographic locality, i.e., they can communicate with datacenters that are 315 faster or have less working load [6, 7, 9]. Hence, system developers have to consider 316 the trade-off between these features for the placement strategy of data replicas across 317 multiple datacenters when they are building geographically distributed services for 318 applications. 319



Fig. 6 The decentralized method to place data replicas across multiple datacenters

### 320 4 Data Striping Methods

### 321 4.1 Introduction

The main purpose of a storage system is to make data persistent, so reliability and availability should be top priority concern for storage systems. Actually, there are a variety of factors that may cause storage system unavailable. For example, if a server fails, the storage system is unable to provide storage services. Some physical damage to a hard disk will result in the loss of data stored. Therefore, it is indispensable for storage systems to introduce some techniques to make them reliable.

A lot of research work has been done in recent years to improve the availabil-328 ity and reliability of storage systems. The main idea is to generate some redundant 329 information of every data block and distribute them on different machines. When 330 one server becomes outage, another server that holds the redundant data can replace 331 the role of the broken server. During this time, the storage system can still provide 332 storage service. When one data block is broken, then other redundant data blocks 333 will restore the broken one. Thus, the availability and reliability is improved. Gen-334 erally, redundant data can be presented in two ways: one is using full data backup 335 mechanism, called full replication; the other is erasure code. 336

Full replication, also known as multi-copy method, is to store multiple replicates 337 338 of data on separate disks, in order to make the data redundant. This method does not involve specialized encoding and decoding algorithms, and it has better fault-339 tolerance performance. But full replication has lower storage efficiency. Storage 340 efficiency is the sum of effective capacity and free capacity divided by raw capacity. 341 When storing N copies of replica, the disk utilization is only 1/N. For relatively 342 large storage systems, full replication brings extra storage overhead, resulting in 343 high storage cost. 344

Along with the increase of the data that a storage system holds, a full replication
 method has been difficult to adapt to mass storage system for redundant mechanism in
 disk utilization and fault tolerance requirements. Therefore, erasure code is becoming
 a better solution for mass storage.

### 349 4.2 Erasure Code Types

Erasure code is derived from communication field. At first, it is mainly used to solve 350 error detection and correction problems in data transmission. Afterwards, erasure 351 code gradually applied to improve the reliability of storage systems. Thus, erasure 352 code has been improved and promoted according to the characters of storage system. 353 The main idea of erasure code is that the original data can be divided into k data 354 fragments, and according to the k data fragments, m redundant fragments can be 355 computed according some coding theory. The original data can be reconstructed by 356 any of the m + k fragments. There are many advantages of erasure code, the foremost 357 of these is the high storage efficiency compared with the mirroring method. 358

There are many types of erasure code. Reed-Solomon code [52] is an MDS code 359 that can meet any number of data disks and redundant disk number. MDS code 360 (maximum distance separable code) is a kind of code that can achieve the theoretically 361 optimal storage utilization. The main idea of Reed Solomon code is to visualize 362 the data encoded as a polynomial. Symbols in data are viewed as coefficients of a 363 polynomial over a finite field. Reed Solomon code is a type of horizontal codes. 364 Horizontal code has the property that data fragments and redundant fragments are 365 stored separately. That is to say, each stripe is neither data stripe nor redundant stripe. 366 Reed Solomon codes are usually divided into two categories: one is Vandermonde 367 Reed Solomon code, and the other is Cauchy Reed Solomon code [53]. The difference 368 between these two categories of Reed Solomon codes is that they are using different 369 generation matrix. For Vandermonde Reed Solomon code, the generation matrix is 370 Vandermonde matrix, and multiplication on Galois filed is needed which is very 371 complex. For Cauchy Reed Solomon code, it is Cauchy matrix, and every operation 372 is XOR operation, which is coding efficient. Figure 7 shows the Encoding principle 373 for Reed-Solomon codes. 374

Compared with Reed Solomon Codes, Array Code [54] is totally based on XOR operation. Due to the efficient of encoding, decoding, updating and reconstruction,



Author's Proof



Fig. 7 Reed-Solomon codes

Array code is widely used. Array code can be categorized as two types due to the placement of data fragment and redundant fragment.

Horizontal parity array codes make data fragments and redundant fragments stored
on different disks. By doing this, Horizontal parity array codes have better scalability.
But most of it can just hold 2 disk failures. It has a drawback on updating data. Every
time updating one data block will result in at least one read and one write operation
on redundant disk. EVENODD code [55] is one kind of Horizontal parity array codes
that used widely.

Vertical parity array codes make data fragment and redundant fragment stored in the same stripe. Because of this design, the efficiency of data update operation will be improved. However, the balance of vertical parity array code leads to a strong interdependency between the disks, which also led to its poor scalability. XCODE [56] is a kind of vertical parity array code, which has theoretically optimum efficiency on data update and reconstruction operation.

### 391 4.3 Erasure Codes in Data Centers

In traditional storage systems such as early GFS and Windows Azure Storage, to ensure the reliability, triplication has been favored because of its ease of implementation. But triplication makes the stored data triple, and storage overhead is a major concern. So many system designers are considering erasure coding as an alternative. Most distributed file systems (GFS, HDFS, Windows Azure) create an append-only write workload for large block size. So data update performance is not a concern.

Using erasure code in distributed file systems, data reconstruction is a major concern. For one data of k data fragment and m redundant fragment, when any one



Fig. 8 LRC codes

of that fragment is broken or lost, to repair that broken fragment, *k* fragments size
of network bandwidth will be needed. So some researchers found that the traditional
erasure code does not fit distributed file system very well. In order to improve the
performance of data repair, there are two ways.

One is reading from fewer fragments. In Windows Azure Storage System, a new set 404 of code called Local Reconstruction Codes (LRC) [57] is adopted. The main idea of 405 LRC is to reduce the number of fragments required to reconstruct the unavailable data 406 fragment. To reduce the number of fragments needed, LRC introduced local parity 407 and global parity. As Fig. 8 shows below,  $x_0$ ,  $x_1$ ,  $x_2$ ,  $y_0$ ,  $y_1$ ,  $y_2$  are data fragments,  $p_x$ 408 is the parity fragment of  $x_0$ ,  $x_1$ ,  $x_2$ .  $p_y$  is the parity fragment of  $y_0$ ,  $y_1$ ,  $y_2$ .  $p_0$  and  $p_1$  are 409 global parity fragments.  $p_x$  and  $p_y$  are called local parity.  $p_0$  and  $p_1$  are global parity. 410 When reconstructing  $x_0$ , instead of reading  $p_0$  or  $p_1$  and other 5 data fragment, it is 411 more efficient to read  $x_1$ ,  $x_2$  and  $p_x$  to compute  $x_0$ . As we can see LRC is not a MDS, 412 but it can greatly reduce the cost of data reconstruction. 413

Another way to improve reconstruction performance is to read more fragments but less data size from each. Regenerating codes [58] provide optimal recovery bandwidth among storage nodes. When reconstructing fragments, it does not just transmit the existing fragments, but sends a liner combination of fragments. By doing this, the recovery data size to send will be reduced. Rotated Reed-Solomon codes [59] and RDOR [60] improve reconstruction performance in a similar way.

### 420 **5** Consistency Models

### 421 5.1 Introduction

422 Constructing a globally distributed system requires many trade-offs between avail 423 ability, consistency, and scalability. Cloud storages are designed to serve for a large
 424 amount of internet-scale applications and platforms simultaneously, which is often

named as infrastructure service. To meet operational requirements, cloud storage
must be designed and implemented as highly available and scalable, in order to
serve consumers requests from all over the world.

One of the key challenges in build cloud storage is to provide a consistency 428 guarantee to all client requests [63]. Cloud storage is a large distributed system 429 deployed world-widely. It has to process millions of requests every hour. All the 430 low-probability accidents in normal systems are often to happen in the datacenters of 431 cloud storage. So all these problems must be taken care of in the design of the system. 432 To guarantee consistent performance and high availability, replication techniques are 433 often used in cloud storage. Although replication solves many problems, it has its 434 costs. Different client requests may see inconsistent states of many replicas. To solve 435 this problem, cloud storage must define a consistency model that all requests to 436 replicas of the same data must follow. 437

Like many widespread distributed systems, cloud storage such as Amazon S3 often provides a weak consistency model called eventual consistency. Different clients may see different orders of updates to the same data object. Some cloud storage like Windows Azure also provides strong consistency that guarantees linearizability of every update from different clients. Details will be discussed in the following sub-sections.

### 444 5.2 Strong Consistency

445 Strong consistency is the most programmer-friendly consistency model. When a
446 client commits an update, every other client would see the update in subsequent
447 operations. Strong consistency can help achieve transparency of a distributed system.
448 When developer uses a storage system with strong consistency, it appears like the
449 system is a single component instead of many collaborating sub-components mixed
450 together.

However, this approach has been proved as difficult to achieve since the middle
of last century, in the database area for the first time. Databases are also systems with
heavy use of data replications. Many of such database systems were design to shut
down completely when it cannot satisfy this consistency because of node failures.
But this is not acceptable for cloud systems, which is so large that small failures are
happening every minute.

457 Strong consistency has its weak points, one of which is that it lowers system 458 availability. In the end of last century, with large-scale Internet systems growing 459 up, designs of consistency model are rethought. Engineers and researchers began to 460 reconsider the tradeoff between system availability and data consistency. In the year 461 of 2000, CAP theorem was introduced [61]. The theorem states that for three prop-462 erties of shared-data systems—data consistency, system availability, and tolerance 463 to network partition—only two can be achieved at any given time.

It is worth noting, that the concept of consistency in cloud storage is different to that in transactional storage systems such as databases. The common ACID property

### 469 5.3 Weak Consistency

According the CAP theory, a system can achieve both consistency and availability, if 470 it does not tolerate network partitions. There many techniques which make this work, 471 one of which is to use transaction protocols like two phase commit. The condition 472 for this is that both client and server of the storage systems must be in the same 473 administrative environment. If partition happens and client cannot observe this, the 474 transaction protocol would fail. However, network partitions are very common in 475 large distributed systems, and as the system scale goes up, the chances of network 476 partition would increase. This is one reason why one cannot achieve consistency and 477 availability at the same time. The CAP theory provides two choices for developers: (1) 478 sticking to strong consistency and allowing system goes unavailable under partitions 479 (2) using relaxed consistency [65] so that system is still available under network 480 partitions. 481

No matter what kind of consistency model the system uses, it requires that ap-482 plication developers are fully aware of the consistency model. Strong consistency is 483 usually the easiest option for client developer. The only problem the developers have 484 to deal with is to tolerate the unavailable situation that might happen to the system. 485 If the system takes relaxed consistency and offers high availability, it may always 486 accept client requests, but client developers have to remember that a write may get 487 its delays and a read may not return the newest write. Then developers have to write 488 the application in a way so that it can tolerant the delay update and stale read. There 489 are many applications that can be design compatible for such relaxed consistency 490 model and work fine. 491

There are two ways of looking at consistency. One is from the developer/client point of view: how they observe data updates. The other is from the server side: how updates flow through the system and what guarantees systems can give with respect to updates.

Let's show consistency models using examples. Suppose we have a storage system 496 which we treat as a black box. To judge its consistency model we have several clients 497 issuing requests to the system. Assume they are client A, client B, client C. All 498 three clients issue both read and write requests to the system. The three clients are 499 independent and irrelevant. They could run on different machines, processes, or 500 threads. The consistency model of the system can be defined by how and when 501 observers (in this case the clients A, B, or C) see updates made to a data object in the 502 storage systems. Assume client A has made an update to a data object: 503

• Strong consistency. After the update completes, any subsequent access (from any of A, B, or C) will return the updated value.

• Weak consistency. The system does not guarantee that subsequent accesses will return the updated value. A number of conditions need to be met before the value will be returned. The period between the update and the moment when it is guaranteed that any observer will always see the updated value is dubbed the inconsistency window.

There are many kinds of weak consistency; we list some of the most common ones as below.

Causal consistency [66]. If client A has communicated to client B that it has updated a data item, a subsequent access by client B will return the updated value, and a write is guaranteed to supersede the earlier write. Access by client C that has no causal relationship to client A is subject to the normal eventual consistency rules.

Eventual consistency [62]. This is a specific form of weak consistency; the storage 518 system guarantees that if no new updates are made to the object, eventually all 519 accesses will return the last updated value. If no failures occur, the maximum size 520 of the inconsistency window can be determined based on factors such as commu-521 nication delays, the load on the system, and the number of replicas involved in the 522 replication scheme. The most popular system that implements eventual consis-523 tency is the domain name system. Updates to a name are distributed according to 524 a configured pattern and in combination with time-controlled caches; eventually, 525 all clients will see the update [64]. 526

- Read-your-writes consistency. This is an important model where client A, after having updated a data item, always accesses the updated value and never sees an older value. This is a special case of the causal consistency model.
- Session consistency. This is a practical version of the previous model, where a client accesses the storage system in the context of a session. As long as the session exists, the system guarantees read-your-writes consistency. If the session terminates because of a certain failure scenario, a new session must be created and the guarantees do not overlap the sessions.
- Monotonic read consistency. If a client has seen a particular value for the object, any subsequent accesses will never return any previous values.

• Monotonic write consistency. In this case, the system guarantees to serialize the writes by the same client. Systems that do not guarantee this level of consistency are notoriously difficult to program.

These consistency models are not exclusive and independent. Some of the above 540 can be combined together. For example, the monotonic read consistency can be 541 combined with session-level consistency. The combination of the both consistencies 542 is very practical for developers in a cloud storage system with eventual consistency. 543 These two properties make it much easier for application developers to build up their 544 apps. They also allow the storage system to keep a relax consistency and provide 545 high availability. As you can see from these consistency models, quite a few different 546 circumstances are possible. Applications need to choose whether or not one can deal 547 with the consequences of particular consistency. 548

### 549 6 Cloud of Multiple Clouds

### 550 6.1 Introduction

Although cloud storage providers claim that their products are cost saving, trouble-551 free, worldwide 24/7 available and reliable, reality shows that (1) such services are 552 sometimes not available to all customers; and (2) customers may experience vastly 553 different accessibility patterns from different geographical locations. Furthermore, 554 there is also a small chance that clients may not even be able to retrieve their data 555 from a cloud provider at all, which usually occurs due to network partitioning and/or 556 temporary failure of cloud provider. For example, authors of [67] reported that this 557 may also cause major cloud service providers to fail providing services for hours or 558 days sometimes. Although cloud providers sign Service Level Agreements (SLA) 559 with their clients to ensure availability of their services, users have complained 560 that these SLAs are sometimes too tricky to break. Moreover, even when a SLA is 561 violated, the compensation is only a minor discount for the payment and not to cover 562 a customer's loss resulted by the violated SLA. 563

Global access experience can be considered as one specifically important issue of 564 availability. In current major cloud storages, users are asked to create region-specific 565 accounts/containers before putting their data blobs/objects into them. The storage 566 provider then stores data blobs/objects into a datacenter in the selected locations; 567 some providers may also create cross-region replicas solely for backup and disaster 568 recovery. A typical result of such topology is an observation where users may experi-569 ence vastly different services based on the network condition between clients and the 570 datacenter holding their required data. Data loss and/or corruption are other impor-571 tant potential threats to users' data should it be stored on a single cloud provider only. 572 Although users of major cloud storage providers have rarely reported data loss and/or 573 corruption, prevention of such problems is not 100 % guaranteed either. Medium to 574 small sized cloud providers may provide a more volatile situation to their customers 575 as they are also in danger of bankruptcy as well. 576

In this section, we present a system named  $\mu$ LibCloud to address the two aforementioned problems of cloud customers; i.e., (1) availability of data as a whole and (2) different quality of services for different customers accessing data from different locations on the globe.  $\mu$ LibCloud is designed and implemented to automatically and transparently stripe data into multiple clouds—similar to RAID's principle in storing local data.  $\mu$ LibCloud is developed based on Apache libCloud project [3], and evaluated through global-wide experiments.

Our main contributions include: (1) to conduct global-wide experiments to show how several possible factors may affect availability and/or global accessibility of cloud storage services; (2) to use erasure codes based on observations. We then design and implement  $\mu$ LibCloud using erasure code to run benchmarks accessing several commercial clouds from different places in the world. The system proved the effectiveness of our method. Author's Proof

Fig. 9 Layer abstraction of cloud storage

Applications
Library
REST/SOAP
Cloud Storage

### 590 6.2 Architecture

Using a "cloud-of-cloud" rationale [68],  $\mu$ LibCloud is to improve availability and 591 global access experience of data. Here the first challenge is how to efficiently and 592 simultaneously use multiple cloud services. They follow different concepts and offer 593 different ways to access their services. As shown in Fig. 9, cloud storage providers 594 usually provide REST/SOAP web service interface to developers along with their 595 libraries for different programming languages for developers to further facilitate 596 building cloud applications. To concurrently use multiple cloud storages, two op-597 tions are available. The first option is to set up proxy among cloud storages and 598 applications. In this case, all data requests need to go through this proxy. To store 599 data in cloud storages, this proxy receives original data from client, divides the data 600 into several shares, and sends each share to different clouds using different libraries. 601 To retrieve data, it fetches data shares from each cloud, rebuilds the data, and sends 602 it back to clients. The second option-more complicated-is to integrate the sup-603 port for multiple cloud storages directly into a new client library-replacing original 604 ones. In this case, client applications only use this newly provided library to connect 605 to all clouds. The main difference between these two options is the transparency in 606 the second option to spread/collect data to/from multiple clouds. 607

The first choice is more straightforward in design; it uses a single layer for extra 608 work, keeps the client neat and clean, includes many original libraries when imple-609 mented, and is usually run on independent servers. It also brings more complexity to 610 system developers to maintain extra servers and their proper functioning. The second 611 choice, on the other hand, benefits developers by providing them a unique tool; this 612 approach also reduces security risk because developers do not need to put their secret 613 keys on the proxy. It however also leads to other challenges on how to design and 614 implement such systems; e.g., how multiple clients can coordinate with each other 615 without extra servers. Furthermore, the client library must be efficient and resource 616 saving because it needs to be run along with application codes. 617

In the design of  $\mu$ LibCloud, we chose to practice the second option so that it will have lesser of a burden on application developers. We also assume that consumers who choose to use cloud storage rather than to build their own infrastructure would not want to set up another server to make everything work. Figure 10 shows the basic architecture of  $\mu$ LibCloud with a single client; this figure also shows how  $\mu$ LibCloud serves upper-level users, while hiding most of development complexities of such systems.



### 625 6.3 Data Striping

As described before, data is first encoded into several original and redundant shares,
 and then stored on different providers. Through this redundancy, data not only is
 protected against possible failures of particular providers—high availability, but also
 tolerates the instability of individual clouds and provides consistent performance.

Among many possible choices for data encoding [69], we choose the most widely used erasure code [70] that is widely used in both storage hardware [71] and distributed systems [72]. Here, coding efficiency is a major concern because all the data striping algorithm work is performed at clients' side; i.e., large overheads that could decrease performance of applications is strongly unacceptable.

Figure 11 shows principles of erasure coding. As can be seen, data is first divided into k equal-sized fragments called original data shares. Then, r parity fragments with the same size as original data shares are computed and called redundant data



Fig. 12 Data stripes stored in each cloud

22

shares. This will generate a total of m = k + r equal-sized shares. The erasure code algorithm guarantees any arbitrary k shares—out of total m shares—is sufficient enough to reconstruct the original data. Both k and r are positive values and are predefined by each user.

Here, we also define redundancy rate as R = m/k to reflect the amount of storage overhead for storing data. For example, if k = 1, m = 2; then, R = m/k = 200 %. It means that each data takes twice of its original size when stored: one original and one replica. In this case, each piece is enough to reconstruct the original data—like RAID 1 (mirroring). If k = 4, m = 5 (like RAID 5); then, R = m/k = 120 %. It means that we need extra 20 % of storage to store any data. In this case, every four pieces—out of all available five pieces—are enough to reconstruct the original data.

In practice, we do not simply just divide an object into several parts and encode them, but the original data is first divided into several chunks, and then erasure coding is performed on each chunk; default chunk's size is usually 64 KB (Fig. 12). There are two benefits in splitting data into several chunks: (1) computation of erasure coding can be parallelized, and (2) reading and writing of file data—such as video and audio—can also be easily supported.

### 655 6.4 Retrieving Strategy

<sup>656</sup> If a developer divides data into (k, m) shares, among all m parts of data the client <sup>657</sup> library only needs k parts to reconstruct the original data. Although retrieving all parts <sup>658</sup> of data could avoid the potential risks of failures, it is unnecessary in most cases. <sup>659</sup> It also wastes more bandwidth and costs more money. Here, although retrieving k<sup>660</sup> data shares to recover the data is enough, selecting the best possible k shares can be <sup>661</sup> tricky. In µLibCloud we offer the following three data fetching strategies.

1. Efficient: Users want to use the *k* most available clouds to retrieve data pieces. Here, to determine which ones are faster,  $\mu$ LibCloud dynamically measures their download speed. When retrieving an object, all metadata files are downloaded first and their link speed is recorded. Upon that, *k* fastest clouds to fetch data are selected. During downloading the main data,  $\mu$ LibCloud keeps recording the download speed to compute its average. The larger data is, the more accurate network estimation would be.

Economical: If application is mainly run in the background—like a backup program storing data into clouds [73]—, users can tolerate spending more time. In such cases, economical cost is more important than speed. μLibCloud also offers a cost-saving mode, in which it will select k providers with lowest prices.

Custom: We also offer an option, allowing developers to set priorities of their
own. This may be preferable in case that they are using computing and storage
resources provided by the same provider. For example, if a developer is deploying
applications into EC2 and use storage of S3, it would be reasonable that s/he wants
to use S3 as the first choice.

### 678 6.5 Mutual Exclusion

When there is more than one client in the system, they must be able to coordinate 679 with each in certain ways to avoid conflicts. Such conflicts can result in not only 680 client read failures, but also inconsistent states and/or even data loss. For example, 681 if two clients concurrently write to the same data file without any locking, they may 682 write to each other's share and produce problems. In the worst case scenario, if the 683 provider takes an eventual consistency model (like Amazon S3), all unordered writes 684 would succeed although only the later ones become effective. As a result, it would be 685 very probable that a client succeeds modifying several data shares, while the other 686 client succeeds in the rest of data shares; both clients would return successful, while 687 data inconsistency has already occurred! The following options are among the most 688 suitable one for our needs. 689

Setting up a central lock server such as ZooKeeper [74] to coordinate all writes.
 This approach is easy and correct for a system like μLibCloud, yet with certain
 flaws. Firstly, with this approach clients need to maintain another system, which
 violates goals and principles of using clouds for simplicity in the first place.

695

696

697

715

24

Secondly, coordinators like ZooKeeper usually have throughput issues because of their leader-follower architecture, especially in internet-scale situation. Although this can be reduced by manually partitioning data onto multiple groups of ZooKeeper systems, this would still make the system extremely complex.

 Running a client-client agreement protocol. Here, instead of deploying an additional central lock service, agreement protocols such as Paxos [75] handles the situation. This approach eliminates the trouble of bringing a lock service, but requires clients to be able to communicate with each other. In this case, frequent membership changes can seriously damage system performance. In fact, this approach is almost the same—in logic—as the first option if each client runs with a ZooKeeper member deployed to the same machine.

3. Manipulating lock files on each cloud storage. Instead of setting up an additional lock server or running an agreement protocol among clients, there is another approach more suitable to this situation. Each client creates empty files on each cloud as lock-files; this is called mutual exclusion in the area of distributed algorithms [76]. This option is more difficult to achieve because each cloud is purely an object storage that offers neither computing ability, nor a common compare and swap (cas) semantics usually used in fulfilling lock services.

<sup>712</sup> In order to achieve mutual exclusion without introducing new bottlenecks, we intro-

duced Algorithm 1 based on the third option. This algorithm is an improved version

of another algorithm formerly designed by Bessani [77].

```
Algorithm 1 Mutual exclusion at client side
  input: id, and providers [1, 2, ..., m]
  output: success or failure
  for i from 1 to m do
    create a file named lock_id on provider[i]
  end for
  count = 0
  for i from 1 to m do
    list all lock files on provider[i]
    if found any lock files created by other client then
      count = count + 1
    end if
  end for
  if count \ge m/2 then
    for i from 1 to m do
      delete lock_id from provider[k]
    end for
  else
    {critical section}
                             // lock succeeds
  end if
```

7181. The algorithm is fault tolerate to possible failures of less than m/2 providers. In719case a client fails and stops during any step, we add a timestamp  $t_{create}$  to the name720of each lock file. Thus, when a client lists a file name with the  $t_{create} + t_{delta} < t_{now}$ ,721s/he can confidently deletes the expiring lock. To maintain correctness, we must722choose a  $t_{delta}$  large enough to cover the entire operation time when created; it723must also be able to tolerate possible time differences among clients.

2. To be correct, the algorithm requires each cloud to have an appropriate consistency 724 model. To be specific, after each 'create' command all lists must see the creation. 725 However, several major cloud providers, such as Amazon S3, employ an eventual 726 consistency model [78]. It means the writes are not visible to reads immediately, 727 and if one client detects a change, it does not imply other clients can also detect 728 it. To tolerate eventual consistency, the client may need to wait for another time 729 period, after each write to make sure it can be seen by all clients too. The time 730 period is set by observation to model time delays among clients [64]. 731

Amazon S3 recently releases an enhanced consistency model to most of its cloud
storages, namely "read-after-write" consistency to ensure that for newly created objects, the write (not overwrite) can be seen immediately. Our algorithm (Algorithm 1)
employs this feature in its locking system; this is why Algorithm 1 creates new lock
files instead of writing to the old ones.

737 3. The algorithm is obstruction-free [79]; i.e., it is still possible—although very
rare—that no client can progress. This flaw could be tolerated because most
rapplications tend to have many more reads than writes—where only very few
writes require mutual exclusion.

## 741 7 Privacy and Security of Storage System

### 742 7.1 Introduction

In the last few years, cloud computing has enabled more and more customers (such as 743 companies or developers) to run their applications on the remote servers with elastic 744 storage capacity and computing resources required on demand. The proliferation 745 of cloud computing encourages customers to store and keep their data in the cloud 746 instead of maintaining local data storage [32-34, 38, 39]. However, a key factor 747 that may hinder the process of data migration from local storage to the cloud is the 748 potential privacy and security concerns inside clouds [33, 34]. As customers don't 749 own and manage remote servers directly by themselves, any malicious applications or 750 administrators in the cloud can get access to, abuse or even damage the data of normal 751 customers' applications. This phenomenon is especially adverse to the confidentiality 752 of sensitive data objects of customers such as banks or financial companies. Under 753 this circumstance, datacenters in the cloud must maintain strong protections on the 754

717

privacy and security of data objects against untrustworthy applications, servers and
 administrators during the process of data storing and accessing [34, 36, 39, 46].

To guarantee data privacy and security in storage system of datacenters in the 757 cloud, several basic solutions such as data access control [38–41], data isolation 758 [36, 37, 42, 46, 47] and cryptographic techniques [35, 40, 43-45] have been proposed 759 by researchers. All these solutions are intended to meet different requirements of 760 data privacy and security and to make even the most privacy and security demanding 761 applications to migrate their sensitive data into cloud with no concerns. In this section, 762 combined with our experience of building privacy and security policies in datacenters 763 in the cloud, we will present how these mechanisms can be used in a real world. 764

### 765 7.2 Fine-Grained Data Access Control

Data access control is highly related to the privacy and security provided to applica-766 tions when they are accessing the data [33, 38, 39, 41]. Applications, if not allowed, 767 don't have the authority to access the data of others. Besides, each application may 768 have its own access control policies to maintain the data privacy and security among 769 its users. For instance, one application may require that only its administrators can 770 have the authority to modify and delete its data and other common users can only read 771 these data. Therefore, storage systems in datacenters must ensure strict and flexible 772 data access control mechanisms for upper applications to secure the data object sets 773 of every application. 774

Algorithm 2 The procedure of storing a data object.
Input: $P_n, d_m$
1. Get $k_m$ of $d_m$
2. Get $L_n$ of $P_n$
3. Insert $k_m$ into $L_n$
4. Keep $d_m$ in physical storage

775 776

Figure 13 illustrates the overview of a fine-grained access control mechanism on 777 data object level in a datacenter. As presented in Fig. 13, there are two main data 778 structures for the correct process of fine-grained data-object-level access control: a 779 set of lists keeping the keys of data objects that belong to each application and a set 780 of tables recording each application's access control policy. Every application owns 781 its list of keys and access control policy table. When one application stores a data 782 object into the datacenter, the storage system will allocate a globally unique key to 783 this data object and add this key into the list of this application, which means this 784 data object does belong to such application. Denote an application as  $P_n$ , the list of 785  $P_n$  as  $L_n$ , a data object as  $d_m$  and the key of data object  $d_m$  as  $k_m$ , then the process of 786 storing data in this mechanism could be summarized as Algorithm 2. 787



Fig. 13 The overview of the fine-grained access control mechanism on the data object level

When an access request for a data object issued from an application arrives at the 788 datacenter, the storage system will first get the key of the data object and verifies 789 if this key is in this application's list. Storage system will forbid the application 790 to access this data object if the verification fails. This procedure ensures that data 791 objects of one application are isolated from the other applications against illegal 792 intrusion. Moreover, if the verification passes, the system will further check if this 793 access request meets the requirements listed in the access control policies table of the 794 application. This will prevent unauthorized application users from abusing operations 795 on data of this application that may potentially damage these data. Applications can 796 set and modify their access control policies according to their own demands and the 797 policy information are recorded in their access control policy tables respectively. The 798 access request is accepted and processed only after the check in the access control 799 policy table successes. Denote an access request as  $R_p$  and access control policies 800 table of application  $P_n$  as  $T_n$ , then the procedure to process an access request can be 801 illustrate as Algorithm 3. 802

With Algorithm 2 and Algorithm 3, the data privacy and security could be achieved 803 across applications through fine-grained data-object-level access control mechanism 804 without impacting the normal usage of data by authorized users of each application. 805 Furthermore, as these two data structures (lists and tables) that are used by the 806 access control mechanism could keep a consistent view across multiple datacenters 807 using replication strategy presented in Sect. 3, the privacy and security of data could 808 be easily guaranteed through this fine-grained data-object-level data access control 809 mechanism among multiple datacenters in the cloud. 810

Algorithm 3 The procedure of processing an access request on a data object.
Input: $P_n, R_p, d_m$
1. Get $k_m$ of $d_m$
2. Get $L_n$ of $P_n$
3. if $k_m \in L_n$ then
4. Get $T_n$ of $P_n$
5. if $R_p$ meets requirements in $T_n$ then
6. Process $R_p$
7. else
8. Refuse to process $R_p$
9. end if
10. else
11. Refuse to process $R_p$
12. end if

811 812

### 813 7.3 Security on Storage Server

814 Under fine-grained data-object-level access control mechanism, the privacy and security of applications' data could be protected against external untrusted users and 815 applications. However, data stored in the storages servers of datacenters are still prone 816 to abuse or compromise by untrusted processes running in these servers or malicious 817 administrators of datacenters that can get the whole authority of the OS [36, 37, 51]. 818 819 To address this issue, most studies [36, 37, 42, 46, 47, 51] use virtual-machine-based protection mechanisms to isolate applications' data kept in hardware (memories and 820 disks) of storage servers from operating systems and other processes, and to authen-821 ticate the integrity of these data. This protection ensures that even operating systems 822 carry out the overall task of managing data they cannot read or modify them. With 823 this guarantee, even though malicious administrators or untrusted processes get the 824 authority of OS, they have no access to abusing or damaging the data stored in the 825 hardware. When trusted applications request to get their data, this mechanism would 826 make sure that these applications will be presented with a normal view of their orig-827 inal data, hiding the complex underlying details of protection. Hence, the privacy 828 and security of applications' data can be maintained in storage servers of datacenters 829 in the cloud. 830

Figure 14 characterizes the architecture of the privacy and security protection 831 mechanism in storage servers. The key component, as shown in Fig. 14, to protect 832 the privacy and security of applications' data in hardware is the virtual machine mon-833 itor (VMM). The VMM could monitor the process/OS interactions such as system 834 calls [36, 42] and directly manage the hardware to isolate memories and disks from 835 operating systems [37, 42, 46], which makes it possible to prevent the data privacy 836 and security against malicious processes or administrators that can get the authority 837 to control operating systems. 838



Fig. 14 The architecture of the privacy and security protection mechanism in a storage server

Generally, each process owns its independent virtual memory address space and 839 is associated with a page table that maps the virtual memory address into the physical 840 memory address [48] to use memory. The page tables of processes and the operations 841 of address mappings are managed by the OS and thus, it has the authority to access 842 the memory address space of all processes running on it. As applications' requests 843 are served by specific processes in storage servers of datacenters, once malicious 844 processes or administrators steal the operating system's authority, they can easily 845 access the data of other normal processes through their page tables and threaten the 846 847 privacy and security of applications' data. To address this challenge, VMM could protect the page tables of each process and complete the operations of memory 848 address mappings instead of operating system [48, 51]. The OS can only access its 849 kernel memory space through its own page table, without interleaving with other 850 processes. However, even though the OS doesn't know the distribution of processes' 851 virtual memory in the physical memory, malicious processes or administrators could 852 also access the physical memory through OS [48, 49] and analyze or tamper the data 853 in the memory [42, 50]. As a result, VMM is responsible for keeping the data in the 854 physical memory in an encrypted and integrated view [49]. When a process requests 855 to put data into memory, VMM will detect this request, encrypt the data and then put 856 857 the encrypted data into the memory. If one trusted process requests to get its data in the memory, VMM will first authenticate the integrity of the encrypted data and then 858 decrypt them before returning the original clear data to this process, which doesn't 859 have to cover this middle process and just utilizes memory as normal. To complete 860 the encrypting and decrypting procedures mentioned above, VMM holds a specific 861

zone of memory that is secure enough against the attacks from operating systems
and processes. Consequently, when processes are serving applications' requests in
memory of storage servers, the privacy and security of their data in memory can be
strongly protected.

As most of applications' data will be stored into disks of storage servers in datacen-866 ters, it is also critical to guarantee the privacy and security of data in disks [36, 42, 51]867 not only because untrusted processes and administrators that get the authority of OS 868 can directly access data in disks through I/O operations, but administrators could 869 fetch disks manually. As a result, data in disks must also be stored in an encrypted 870 871 view so that even some processes or administrators get control on the disks of storage servers, they have no way to abuse or compromise the data stored in them. VMM also 872 has the responsibility for data encryption/decryption when processes interact with 873 disks through the OS. When a process wants to write its data into disks, it will use 874 a system call sys write [48] and passes the data to the operating system, which will 875 execute the operations to really write data to disks. VMM will detect this system call 876 from the process and obtain the data before passing to the operating system. Then 877 VMM encrypts these data and calculate the checksum of the encrypted data for future 878 integrity verification. After this procedure, VMM will transfer the encrypted data to 879 the operating system that will normally write these data into disks. Similarly, when 880 881 one process requests to get its data from disks, it will issue a system call sys\_read to the operating system to fetch these data. VMM will also detect this system call 882 and wait for the operating system to complete the read operations of the encrypted 883 data. Then VMM authenticates the integrity of the encrypted data, decrypts them and 884 return plain data to the process. All the underlying details of encryption/decryption 885 are still hidden from the processes and to the operating system, although it manages 886 the data during the operations of read and write, it only views data after encryption 887 and can't threaten the privacy and security of the original data objects. 888

With these virtual-machine-based mechanisms, the data of applications can be kept in storage servers of datacenters without concerns of being abused or compromised by malicious processes or administrators in the datacenters. As data privacy and security can be achieved in hardware of each storage server, datacenters in the cloud have the ability to provide high privacy and security for applications to move their large sets of data into cloud and freely access their data on demands.

### 895 8 Conclusion and Future Directions

In this chapter we mainly discussed the architecture of modern cloud storage and several key techniques used in building such systems. Cloud storage systems are typically large distributed systems composed of thousands of machines and network devices over many datacenters across multiple continents. Cloud storage and cloud computing are the very mixture of modern storage and network technology. To build and maintain such systems calls for large amount of efforts from numerous developers and maintainers. Although we have discussed about replication, data striping, data consistency, security and some other issues, there are still much more of the iceberg we have not touched. Many conventional techniques in traditional storage techniques applied in cloud storage have the potentiality to evolve, such as the example we give about cloud-of-clouds, which arise from the traditional RAID system. To summarize, cloud storage is a valued area in both practice and research, and the goal of this chapter is to provide a glimpse into it when it grows into a global scale.

### 910 **References**

- Varia, Jinesh. "Cloud architectures." White Paper of Amazon, jineshvaria. s3. amazonaws.
   com/public/cloudarchitectures-varia. pdf (2008).
- 2. Brad Calder, Ju Wang, Aaron Ogus, Niranjan Nilakantan, Arild Skjolsvold, Sam McKelvie, 913 Yikang Xu, Shashwat Srivastav, Jiesheng Wu, Huseyin Simitci, Jaidev Haridas, Chakravarthy 914 Uddaraju, Hemal Khatri, Andrew Edwards, Vaman Bedekar, Shane Mainali, Rafay Ab-915 basi, Arpit Agarwal, Mian Fahim ul Haq, Muhammad Ikram ul Haq, Deepali Bhardwaj, 916 Sowmya Dayanand, Anitha Aduzsumilli, Marvin McNett, Sriram Sankaran, Kavitha Mani-917 vannan, Leonidas Rigas. Windows Azure Storage: a highly available cloud storage service 918 with strong consistency. Proceedings of the Twenty-Third ACM Symposium on Operating 919 Systems Principles (SOSP'11), pages 143–157, 2011. 920
- G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Vosshall, W. Vogels. Dynamo: amazon's highly available key-value store.
   Proceedings of twenty-first ACM SIGOPS Symposium on Operating Systems Principles (SOSP' 07), pages 205–220, 2007.
- Sanjay Ghemawat, Howard Gobioff, Shun-Tak Leung. The Google file system. Proceedings of the nineteenth ACM Symposium on Operating Systems Principles (SOSP' 03), pages 29–43, 2003.
- 5. Chang, Fay, et al. "Bigtable: A distributed storage system for structured data." ACM
   Transactions on Computer Systems (TOCS) 26.2 (2008): 4.
- James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, J. J.
   Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, Wilson
   Hsieh, Sebastian Kanthak, Eugene Kogan, Hongyi Li, Alexander Lloyd, Sergey Melnik, David
   Mwaura, David Nagle, Sean Quinlan, Rajesh Rao, Lindsay Rolig, Yasushi Saito, Michal Szy maniak, Christopher Taylor, Ruth Wang, Dale Woodford, D. Woodford. Spanner: Google's
   globally-distributed database. Proceedings of the 10th USENIX conference on Operating
   Systems Design and Implementation (OSDI' 12), pages 251–264, 2012.
- Yair Sovran, Russell Power, Marcos K. Aguilera, Jinyang Li. Transactional storage for georeplicated systems. Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles (SOSP'11), pages 385–400, 2011.
- 8. Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, David G. Andersen. Don't settle for eventual: scalable causal consistency for wide-area storage with COPS. Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles (SOSP'11), pages 401–416, 2011.
- 9. Cheng Li, Daniel Porto, Allen Clement, Johannes Gehrke, Nuno Preguica, Rodrigo Rodrigues. Making geo-replicated systems fast as possible, consistent when necessary. Proceedings of the 10th USENIX conference on Operating Systems Design and Implementation (OSDI'12), pages 265–278, 2012.
- 48 10. Luiz André Barroso, Urs Hölzle. The Datacenter as a Computer: An Introduction
   49 to the Design of Warehouse-Scale Machines. Morgan & Claypool Publishers, DOI:
   40.2200/S00193ED1V01Y200905CAC006, 2009.

- 11. Jason Baker, Chris Bond, James C. Corbett, JJ Furman, Andrey Khorlin, James Larson,
   Jean-Michel Leon, Yawei Li, Alexander Floyd, Vadim Yushprakh. Megastore: Providing Scal able, Highly Available Storage for Interactive Services. In 5th Conference on Innovative Data
   Systems Research, pages 223–234, 2011.
- Avinash Lakshman, Prashant Malik. Cassandra: a decentralized structured storage system.
   ACM SIGOPS Operating Systems Review, 44(2), pages 35–40, 2010.
- D. B. Terry, M. M. Theimer, Karin Petersen, A. J. Demers, M. J. Spreitzer, C. H. Hauser. Managing update conflicts in Bayou, a weakly connected replicated storage system. Proceedings of the fifteenth ACM Symposium on Operating Systems Principles (SOSP'95), pages 172–182, 1995.
- 14. Brian F. Cooper, Raghu Ramakrishnan, Utkarsh Srivastava, Adam Silberstein, Philip
   Bohannon, Hans-Arno, Nick Puz, Daniel Weaver, Ramana Yerneni. PNUTS: Yahoo!'s hosted
   data serving platform. Proceedings of the VLDB Endowment, 1(2), pages 1277–1288, 2008.
- Fred B. Schneider. Implementing fault-tolerant services using the state machine approach: a tutorial. ACM Computing Surveys (CSUR), 22(4), pages 299–319, 1990.
- 16. Leslie Lamport. Paxos made simple. ACM SIGACT News Distributed Computing Column,
   32(4), pages 18–25, 2001.
- Tushar D. Chandra, Robert Griesemer, Joshua Redstone. Paxos made live: an engineering per spective. Proceedings of the twenty-sixth annual ACM Symposium on Principles of Distributed
   Computing, pages 398–407, 2007.
- 18. Mike Burrows. The Chubby lock service for loosely-coupled distributed systems. Proceedings of the 7th symposium on Operating Systems Design and Implementation (OSDI'06), pages 335–350, 2006.
- 974 19. Jeff Dean. Designs, Lessons, and Advice from Building Large Distributed Systems. Keynote
   975 from LADIS, 2009.
- Stacy Patterson, Aaron J. Elmore, Faisal Nawab, Divyakant Agrawal, Amr El Abbadi. Se rializability, not serial: concurrency control and availability in multi-datacenter datastores.
   Proceedings of the VLDB Endowment, 5(11), PAGES 1459–1470, 2012.
- 979 21. Werner Vogels. Eventually consistent. Communications of the ACM—Rural engineering
   980 development, 52(1), pages 40–44, 2009.
- 22. Lisa Glendenning, Ivan Beschastnikh, Arvind Krishnamurthy, Thomas Anderson. Scalable
   consistency in Scatter. Proceedings of the Twenty-Third ACM Symposium on Operating
   Systems Principles (SOSP'11), pages 15–28, 2011.
- Daniel Ford, François Labelle, Florentina I. Popovici, Murray Stokely, Van-Anh Truong, Luiz Barroso, Carrie Grimes, Sean Quinlan. Availability in globally distributed storage systems. Proceedings of the 9th USENIX conference on Operating Systems Design and Implementation (OSDI'10), No. 1–7, 2010.
- Sage A. Weil, Scott A. Brandt, Ethan L. Miller, Carlos Maltzahn. CRUSH: controlled, scalable,
   decentralized placement of replicated data. Proceedings of the 2006 ACM/IEEE conference on
   Supercomputing, 2006.
- Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, David G. Andersen. Stronger Semantics for Low-Latency Geo-Replicated Storage. Proceedings of the 10th USENIX Symposium on Networked Systems Design and Implementation (NSDI'13), 2013.
- Sharad Agarwal, John Dunagan, Navendu Jain, Stefan Saroiu, Alec Wolman, Harbinder
   Bhogan. Volley: automated data placement for geo-distributed cloud services. Proceedings of
   the 7th USENIX conference on Networked Systems Design and Implementation (NSDI'10),
   2010.
- Anton Beloglazov, Rajkumar Buyya. Energy Efficient Resource Management in Virtualized
   Cloud Data Centers. Proceedings of the 2010 10th IEEE/ACM International Conference on
   Cluster, Cloud and Grid Computing, pages 826–831, 2010.
- 28. Zhichao Li, Kevin M. Greenan, Andrew W. Leung, Erez Zadok. Power Consumption in Enterprise-Scale Backup Storage Systems. Proceedings of the Tenth USENIX Conference on File and Storage Technologies (FAST '12), pages 65–71, 2012.

- 29. J. G. Koomey. Growth in data center electricity use 2005 to 2010. Technical report, Standord University, 2011.
- 30. Yi Lin, Bettina Kemm, Marta Patiño-Martínez, Ricardo Jiménez-Peris. Middleware based data
   replication providing snapshot isolation. Proceedings of the 2005 ACM SIGMOD international
   conference on Management of data, pages 419–430, 2005.
- 31. Daniel Peng, Frank Dabek. Large-scale incremental processing using distributed transactions and notifications. Proceedings of the 9th USENIX conference on Operating Systems Design and Implementation, 2010.
- 32. Cong Wang, Qian Wang, and Kui Ren, Wenjing Lou. Privacy-Preserving Public Auditing for
   Data Storage Security in Cloud Computing. Proceedings of IEEE INFOCOM, 2010.
- 33. S. Subashini, V. Kavitha. A survey on security issues in service delivery models of cloud computing. Journal of Network and Computer Applications, 34(1), pages 1–11, 2011.
- 34. H. Takabi, J.B.D. Joshi, G. Ahn. Security and Privacy Challenges in Cloud Computing
   Environments. IEEE Security and Privacy, 8(6), pages 24–31, 2010.
- 35. Kevin D. Bowers, Ari Juels, Alina Oprea. HAIL: a high-availability and integrity layer for
   cloud storage. Proceedings of the 16th ACM conference on Computer and Communications
   Security (CCS'09), pages 187–198, 2009.
- 36. Fengzhe Zhang, Jin Chen, Haibo Chen, Binyu Zang. CloudVisor: retrofitting protection of
   virtual machines in multi-tenant cloud with nested virtualization. Proceedings of the Twenty Third ACM Symposium on Operating Systems Principles (SOSP'11), pages 203–216, 2011.
- 37. Xiaoxin Chen, Tal Garfinkel, E. Christopher Lewis, Pratap Subrahmanyam, Carl A. Wald-spurger, Dan Boneh, Jeffrey Dwoskin, Dan R.K. Ports. Overshadow: a virtualization-based approach to retrofitting protection in commodity operating systems. Proceedings of the 13th international conference on Architectural Support for Programming Languages and Operating Systems, pages 2–13, 2008.
- 38. Wassim Itani, Ayman Kayssi, Ali Chehab. Privacy as a Service: Privacy-Aware Data Storage
   and Processing in Cloud Computing Architectures. Proceedings of Eighth IEEE International
   Conference on Dependable, Autonomic and Secure Computing, pages 711–716, 2009.
- 39. Shucheng Yu, Cong Wang, Kui Ren, Wenjing Lou. Achieving Secure, Scalable, and Fine grained Data Access Control in Cloud Computing. Proceedings of IEEE INFOCOM, 2010.
- 40. Vipul Goyal, Omkant Pandey, Amit Sahai, Brent Waters. Attribute-based encryption for
   fine-grained access control of encrypted data. Proceedings of the 13th ACM conference on
   Computer and Communications Security (CCS'06), pages 89–98, 2006.
- Myong H. Kang, Joon S. Park, Judith N. Froscher. Access control mechanisms for inter organizational workflow. Proceedings of the sixth ACM symposium on Access Control Models
   and Technologies, pages 66–74, 2001.
- H. Chen, F. Zhang, C. Chen, Z. Yang, R. Chen, B. Zang, P. Yew, and W. Mao. Tamper-resistant
   execution in an untrusted operating system using a virtual machine monitor. Parallel Processing
   Institute Technical Report, Number: FDUPPITR-2007-0801, Fudan University, 2007.
- 43. Lein Harn, Hung-Yu Lin. A cryptographic key generation scheme for multilevel data security.
   Computer & Security, 9(6), pages 539–546, 1990.
- Christian Cachin, Klaus Kursawe, Anna Lysyanskaya, Reto Strobl. Asynchronous verifiable se cret sharing and proactive cryptosystems. Proceedings of the 9th ACM conference on Computer
   and Communications Security (CCS'02), pages 88–97, 2002.
- 45. Phillip Rogaway. Bucket hashing and its application to fast message authentication. CRYPTO,
   volume 963 of LNCS, pages 29–42, 1995.
- 46. David Lie, Chandramohan A. Thekkath, Mark Horowitz. Implementing an untrusted operating
   system on trusted hardware. Proceedings of the nineteenth ACM Symposium on Operating
   Systems Principles (SOSP'03), pages 178–192, 2003.
- 47. Stephen T. Jones, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau. Geiger: monitoring
   the buffer cache in a virtual machine environment. Proceedings of the 12th international con ference on Architectural Support for Programming Languages and Operating Systems, pages
   14–24, 2006.

- 48. Abraham Silberschatz, Peter Baer Galvin, Greg Gagne. Operating System Concepts. John
   Wiley & Sons, 2009.
- 49. Guillaume Duc, Ronan Keryell. CryptoPage: an Efficient Secure Architecture with Memory
   Encryption, Integrity and Information Leakage Protection. Proceedings of the 22nd Annual
   Computer Security Applications Conference (ACSAC'06), pages 483–492, 2006.
- David Lie, Chandramohan Thekkath, Mark Mitchell, Patrick Lincoln, Dan Boneh, John
   Mitchell, Mark Horowitz. Architectural support for copy and tamper resistant software. ACM
   SIGPLAN Notices, 35(11), pages 168–177, 2000.
- 1065 51. Hou Qinghua, Wu Yongwei, Zheng Weimin, Yang Guangwen. A Method on Protection of User
   Data Privacy in Cloud Storage Platform. Journal of Computer Research and Development,
   48(7), pages 1146–1154, 2011.
- 1068 52. Reed I S, Solomon G. Polynomial codes over certain finite fields [J]. Journal of the Society for
   1069 Industrial & Applied Mathematics, 1960, 8(2): 300–304.
- 1070 53. Roth R M, Lempel A. On MDS codes via Cauchy matrices [J]. Information Theory, IEEE
   1071 Transactions on, 1989, 35(6): 1314–1319.
- 1072 54. Blaum M, Farrell P, Tilborg H. Array Codes [M]. Amsterdam, Netherlands: Elsevier Science
   1073 B V, 1998.
- 1074 55. Blaum M, Brady J, Bruck J, et al. EVENODD: An efficient scheme for tolerating double disk
   1075 failures in RAID architectures[J]. Computers, IEEE Transactions on, 1995, 44(2): 192–202.
- 56. Xu L, Bruck J. X-code: MDS array codes with optimal encoding[J]. Information Theory, IEEE
   Transactions on, 1999, 45(1): 272–276.
- 1078 57. Huang, Cheng, et al. "Erasure coding in windows azure storage." USENIX ATC. 2012.
- 58. Dimakis A G, Godfrey P B, Wu Y, et al. Network coding for distributed storage systems[J].
  Information Theory, IEEE Transactions on, 2010, 56(9): 4539–4551.
- 59. Khan, Osama, et al. "Rethinking erasure codes for cloud file systems: Minimizing I/O for
   recovery and degraded reads." Proc. of USENIX FAST. 2012.
- Kiang, Liping, et al. "Optimal recovery of single disk failure in RDP code storage systems."
   ACM SIGMETRICS Performance Evaluation Review. Vol. 38. No. 1. ACM, 2010.
- 1085 61. Brewer, Eric A. "Towards robust distributed systems." PODC. 2000.
- 1086 62. Vogels, Werner. "Eventually consistent." Communications of the ACM 52.1 (2009): 40–44.
- 63. Birman, Kenneth P. "Consistency in Distributed Systems." Guide to Reliable Distributed
   Systems. Springer London, 2012. 457–470.
- 64. Bermbach, David, and Stefan Tai. "Eventual consistency: How soon is eventual? An evaluation
   of Amazon S3's consistency behavior." Proceedings of the 6th Workshop on Middleware for
   Service Oriented Computing. ACM, 2011.
- Kanal School, Schol, School, School, School, School, School, School, School, Scho
- 66. Adve, Sarita V., and Kourosh Gharachorloo. "Shared memory consistency models: A tutorial."
   computer 29.12 (1996): 66–76.
- 67. Serious cloud failures and disasters of 2011. http://www.cloudways.com/blog/cloud-failures disastersof-2011/.
- 68. D. Bernstein, E. Ludvigson, K. Sankar, S. Diamond, and M. Morrow, "Blueprint for the intercloud—protocols and formats for cloud computing interoperability," Internet and Web Applications and Services, International Conference on, vol. 0, pp. 328–336, 2009.
- 69. R. G. Dimakis, P. B. Godfrey, Y. Wu, M. O. Wainwright, and K. Ramch, "Network coding for
   distributed storage systems," in In Proc. of IEEE INFOCOM, 2007.
- 70. L. Rizzo, "Effective erasure codes for reliable computer communication protocols," SIG COMM Comput. Commun. Rev., vol. 27, no. 2, pp. 24–36, Apr. 1997. [Online]. Available:
   http://doi.acm.org/10.1145/263876.263881
- 1106 71. H.P.Anvin. The mathematics of raid-6. http://kernel.org/pub/linux/kernel/people/hpa/raid6.pdf.
- 1107 72. H. Weatherspoon and J. Kubiatowicz, "Erasure coding vs. replication: A quantitative
   1108 comparison," Peer-to-Peer Systems, pp. 328–337, 2002.

Author's Proof

- 73. M. Vrable, S. Savage, and G. M. Voelker, "Cumulus: Filesystem backup to the cloud," Trans.
  Storage, vol. 5, no. 4, pp. 14:1–14:28, Dec. 2009. [Online]. Available: http://doi.acm.org/
  10.1145/1629080.1629084
- 74. P. Hunt, M. Konar, F. P. Junqueira, and B. Reed, "Zookeeper: waitfree coordination for internet-scale systems," in Proceedings of the 2010 USENIX conference on USENIX annual technical conference, ser. USENIXATC'10. Berkeley, CA, USA: USENIX Association, 2010, pp. 11–11.
- 1116 75. L. Lamport, "Paxos made simple," ACM SIGACT News, vol. 32, no. 4, pp. 18–25, 2001.
- 1117 76. N. A. Lynch, Distributed algorithms. Morgan Kaufmann, 1996.
- 77. A. Bessani, M. Correia, B. Quaresma, F. Andr'e, and P. Sousa, "Depsky: dependable and secure storage in a cloud-of-clouds," in Proceedings of the sixth conference on Computer systems, ser. EuroSys'11. New York, NY, USA: ACM, 2011, pp. 31–46.
- 1121 78. W. Vogels, "Eventually consistent," Communications of the ACM, vol. 52, no. 1, pp. 40–44, 2009.
- 79. M. Herlihy, V. Luchangco, and M. Moir, "Obstruction-free synchronization: Double-ended queues as an example," in Distributed Computing Systems, 2003. Proceedings. 23rd International Conference on. IEEE, 2003, pp. 522–529.
- 1126 80. Csc cloud usage index. http://www.csc.com/.
- 1127 81. D. Ionescu. (Oct. 2009) Microsoft red-faced after massive sidekick data loss. pcworld.

AQ1. We have inserted city and country in all the author's affiliations. Please check.