

Exploiting Workload Characteristics and Service Diversity to Improve the Availability of Cloud Storage Systems

Bo Mao, *Member, IEEE*, Suzhen Wu, *Member, IEEE*, and Hong Jiang, *Fellow, IEEE*

Abstract—With the increasing utilization and popularity of the cloud infrastructure, more and more data are moved to the cloud storage systems. This makes the availability of cloud storage services critically important, particularly given the fact that outages of cloud storage services have indeed happened from time to time. Thus, solely depending on a single cloud storage provider for storage services can risk violating the service-level agreement (SLA) due to the weakening of service availability. This has led to the notion of Cloud-of-Clouds, where data redundancy is introduced to distribute data among multiple independent cloud storage providers, to address the problem. The key in the effectiveness of the Cloud-of-Clouds approaches lies in how the data redundancy is incorporated and distributed among the clouds. However, the existing Cloud-of-Clouds approaches utilize either replication or erasure codes to redundantly distribute data across multiple clouds, thus incurring either high space or high performance overheads. In this paper, we propose a hybrid redundant data distribution approach, called HyRD, to improve the cloud storage availability in Cloud-of-Clouds by exploiting the workload characteristics and the diversity of cloud providers. In HyRD, large files are distributed in multiple cost-efficient cloud storage providers with erasure-coded data redundancy while small files and file system metadata are replicated on multiple high-performance cloud storage providers. The experiments conducted on our lightweight prototype implementation of HyRD show that HyRD improves the cost efficiency by 33.4 and 20.4 percent, and reduces the access latency by 58.7 and 34.8 percent than the DuraCloud and RACS schemes, respectively.

Index Terms—Cloud-of-clouds, replication, erasure codes, availability, cost efficiency

1 INTRODUCTION

WITH the increasing popularity and cost-effectiveness of cloud storage, many companies and organizations have moved or planned to move data out of their own data centers into the cloud. Typical usage examples include storing backup data and online digital media, such as the recent announcement by the United States Library of Congress to move its digitized content to the cloud [1] and Netflix's dependence on the Amazon S3 storage [2] for the storage of its content. However, solely depending on a particular cloud storage provider has a number of potentially serious problems. First, it can cause the so-called vendor lock-in problem for the customers [3], [4], which results in prohibitively high cost for clients to switch from one provider to another as elaborated in Section 2.1. Second, it can cause service disruptions, which in turn will lead to SLA violation, due to cloud outages, resulting in penalties, monetary or other forms, for the service providers. Examples include a

series of high-profile cloud outages in the year of 2014 for cloud providers, such as Amazon, Microsoft and Google [5], from a 5-minute failure that costs half a million dollars to a week-long disruption that costs an immeasurable amount of brand damage. From January to March 2014, DropBox has experienced two times of service outages. Third, solely depending on a particular cloud storage provider can also result in possible increased service costs and data security issues, such as the data leakage problem [6]. Thus using multiple independent cloud providers, so called Cloud-of-Clouds, is an effective way to provide better availability for the cloud storage systems.

In a Cloud-of-Clouds, data redundancy is introduced to judiciously distribute data among the clouds. Thus, the redundant data distribution scheme is critically important for storage availability, performance, cost and space efficiency. Replication achieves the goals of availability and market mobility, but at a very high storage and bandwidth cost for large files. A more economical approach is to spread the data across multiple providers by introducing erasure-code redundancy to tolerate possible failures or outages, such as RACS [3] and NCCloud [7]. However, these schemes suffer from performance degradation due to the small updates over the networked storage because of the well-known write-amplification problem [8]. For example, a small update in the RACS system will incur a total of four accesses, including traffic of two reads and two writes over the network. Furthermore, a recent study conducted on Facebook's warehouse cluster [9], [10] reveals that more than 180 TB of data is transferred through the top-of-rack

- B. Mao is with the Software School of Xiamen University, Xiamen, Fujian, China. E-mail: maobo@xmu.edu.cn.
- S. Wu is with the Computer Science Department of Xiamen University, Xiamen, Fujian, China. E-mail: suzhen@xmu.edu.cn.
- H. Jiang is with the Department of Computer Science & Engineering at the University of Texas at Arlington, and the Department of Computer Science & Engineering at the University of Nebraska-Lincoln. E-mail: hong.jiang@uta.edu.

Manuscript received 11 June 2015; revised 20 Aug. 2015; accepted 22 Aug. 2015. Date of publication 31 Aug. 2015; date of current version 15 June 2016.

Recommended for acceptance by H. Jin.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.

Digital Object Identifier no. 10.1109/TPDS.2015.2475273

switches everyday for RS-coded data recovery. In other words, the recovery operations consume a large amount of cross-rack bandwidth, thereby rendering the bandwidth unavailable for the foreground jobs. Thus, the data recovery process of the erasure-coded schemes will also incur significant network traffic due to the recovery I/Os in the cloud storage systems.

When examining replication- and erasure-code-based schemes in the context of cloud storage systems, I/O performance and space efficiency are two important metrics. Given the explosive growth in data volume with big data analytics, the I/O bottleneck has become an increasingly daunting challenge in terms of both performance and storage capacity. Besides the normal I/O performance, two additional performance-critical operations emerge in a Cloud-of-Clouds: the degraded read operation due to and the recovery operation from a single cloud service outage. Moreover, Recent IDC studies indicate that in the past five years the volume of data has increased by almost nine times to 7 ZB per year and a more than 44-fold growth to 35 ZB is expected in the next ten years [11]. Managing the data deluge on storage to support (near) real-time data analytics becomes an increasingly critical challenge for Big Data analytics in the Cloud.

On the other hand, previous studies on the workload characteristics have shown that files are of mixed sizes with both small and large files [12], [13]. Moreover, file metadata accesses are much more frequent than file accesses and account for more than half of all the user operations [12], [13]. Thus, the performance of file metadata accesses is critical to the overall system performance and directly affects user experience. The recent studies, including RACS [3], DuraCloud [14], DepSky [4], NCCloud [7], and our own analysis detailed in Section 2 indicate that replication-based schemes are performance-friendly to small files and file metadata while erasure-code-based schemes are performance-friendly and cost-efficient to large files. This suggests that, a sensible data distribution scheme in the Cloud-of-Clouds should dynamically utilize replication and erasure codes based on different file characteristics. However, the existing Cloud-of-Clouds schemes, such as RACS [3], DuraCloud [14], DepSky [4], and NCCloud [7], are purely replication-based or erasure-code-based schemes. Moreover, the diversity characteristics of cloud storage providers and workload characteristics are not fully exploited in their designs.

To address the important storage availability issue in the Cloud-of-Clouds, we propose a hybrid data distribution approach, called HyRD, by considering both the diversity characteristics of cloud storage providers and the workload characteristics. Different from the existing Cloud-of-Clouds approaches, HyRD utilizes replication to store the small files and file system metadata, and erasure codes to store the large files on multiple cloud storage providers. By exploiting the workload characteristics and the diversity of cloud providers, the advantages of both the erasure codes and replication are exploited and their disadvantages are alleviated. The extensive trace-driven experiments conducted on our lightweight prototype implementation of HyRD show that HyRD significantly outperforms RACS and DuraCloud in the I/O performance measure of average response times. Moreover, our evaluation and analysis results also show that HyRD achieves comparable or better cost and space efficiency.

More specifically, this paper makes the following main contributions:

- We investigate the cloud storage diversity from the viewpoint of both the performance and cost characteristics. In addition to the listed prices, we also show the performance results from our own performance evaluations.
- We propose a hybrid data distribution method, HyRD, which incorporates both replication and erasure code schemes. Moreover, HyRD exploits the cloud storage diversity and workload characteristics to distribute data blocks with replication or erasure code schemes judiciously and adaptively.
- We conduct experiments on our lightweight prototype implementation to evaluate the HyRD performance and compare it with the other Cloud-of-Clouds schemes. Moreover, we also provide the cost analysis results based on the prices that are listed in the official websites of the relevant cloud storage providers.

The rest of this paper is organized as follows. Background and Motivation are presented in Section 2. We describe the HyRD architecture and design in Section 3. The performance evaluation is presented in Section 4. We review the related work in Section 5 and conclude this paper in Section 6.

2 BACKGROUND AND MOTIVATION

In this section, we first present some important observations drawn from previous reports about the outages of single Cloud service, the diversity of cloud storage, and the workload characteristics. Then we present the characteristics of the replication- and erasure-code-based redundant data distribution in Cloud-of-Clouds to motivate the HyRD study.

2.1 The Outages of Single Cloud Storage Service

As cloud services reshape the IT landscape, the looming threat of a single cloud storage outage continues to hinder the full embrace of the cloud. When you put your data in the cloud, you lose in control. And the security concerns are considerable. But nowhere is the nightmare as vivid as it is when your cloud storage service goes down. The reason is that cloud storage outages really happen, even for major providers such as Amazon S3, Microsoft Azure, Google, Rackspace and other IaaS vendors. For example, a series of high-profile cloud outages in the year of 2014 happened in major cloud providers, such as Amazon, Microsoft and Google [5], from a 5-minute failure that costs half a million dollars to a week-long disruption that costs an immeasurable amount of brand damage. From January to March 2014, DropBox has experienced two times of service outages [5].

The outages of single cloud storage service directly affect the availability of cloud storage services. The cloud storage outage will lead to possible data loss or unavailability for the users if their cloud storage provider goes out of business or suffers a service outage. Despite of the strict Service-Level Agreements (SLAs) between the cloud provider and the user, service failures and outages do occur and are almost unavoidable. A study conducted by ESG (Enterprise Strategy Group) research has shown that about 58 percent

TABLE 1

Monthly Price Plans (in US Dollars, \$1 = ¥ 6.1) for Amazon S3, Windows Azure Storage, Aliyun Open Storage Service and RackSpace Cloud Files, as of September, 10th 2014 in the China Region

<i>Operations & Vendors</i>	<i>Amazon S3 [2]</i>	<i>Windows Azure [15]</i>	<i>Aliyun [16]</i>	<i>RackSpace [17]</i>
Storage (per GB/month)	\$0.033	\$0.157	\$0.029	\$0.13
Data In (per GB)	Free	Free	Free	Free
Data Out to Internet (per GB)	\$0.201	Free	\$0.123	Free
Put, Copy, Post, and List (per 10 K transactions)	\$0.047	Free	\$0.0016	Free
Get and others (per 10 K transactions)	\$0.0037	Free	\$0.0016	Free

of professionals in SMBs (Small and Medium Businesses) can tolerate no more than four hours of downtime before experiencing significant adverse effect [18], [19]. More seriously, EMC's Disaster Recovery Survey in 2013 [20] has observed that the average cost per hour of downtime is much higher than ever before and 54 percent users suffered from lost data or service downtime, which further stresses the importance of the service/data availability in cloud storage systems.

To address the associated problems induced by a single individual cloud storage provider, a Cloud-of-Clouds solution is proposed in the literature [3], [4], [7], [8], [14]. It redundantly distributes data across multiple providers by means of data redundancy schemes, such as replication and erasure codes. As a result, users can maintain their mobility while insuring against outages of a single individual cloud provider.

2.2 The Diversity of Cloud Storage Services

The services provided by the cloud storage are diverse [3], [21]. The cloud storage providers offer different pricing with different performance characteristics. Table 1 shows the monthly price plans for four major providers as of September 10th 2014. For all the cloud providers, we use the prices from the first chargeable usage tier in the China region (i.e., storage usage within 1 TB/month in Amazon S3; the volume of data transferred out ranges between 1 GB/month and 10 TB/month). We can see that the charged costs of cloud storage providers are different in the aspects of storage, data in/out, and the metadata operations. Moreover, we also conducted the performance evaluations of different cloud storage providers, as shown in Fig. 7 in Section 4. We can see that the access latencies of read and write operations are different for the four cloud storage providers.

Moreover, some cloud storage providers may include extra features such as geographic data distribution, access through mountable file systems, or specific APIs. Changes in these features, or the emergence of new providers with more powerful and attractive characteristics, might compel some users to switch from one provider to another. However, moving from one provider to another one may be very expensive because the switching cost is proportional to the amount of data that has been stored in the original provider [3]. The more data has been stored in the original provider, the higher the switching cost will be paid to the bandwidth cost of data migration. This puts the users at a disadvantage, that is, when the cloud storage provider that has stored the user's data raises prices or negotiates a new contract less favorable to the user, the user has no choice

but to accept because of the high switching cost, hence the so called *vendor lock-in problem* [3], [4].

In addition to the prohibitively high cost for clients to switch from one provider to another, vendor lock-in also subjects users to the possibility of data loss if their cloud storage provider goes out of business or suffers a catastrophe. Recent incidents have shown how failures at cloud storage providers can result in mass data loss for users, and that outages, though seldom, can last up to several hours, even up to days [5].

2.3 The Workload Characteristics

Understanding the workload characteristics is very crucial to avoid any design inefficiencies in a storage system. Previous studies on the workload characteristics have shown two important observations [12], [13], [22]. First, more than 50 percent of files are smaller than 4 KB [12] and only account for less than 20 percent of total storage capacity. The files whose size ranges from 3 to 9 MB accounts for more than 80 percent of the total storage capacity. The larger I/O requests can result in faster throughput by exploiting the access parallelism. The I/O request size is an important factor in workload characterization. Knowing the I/O request size can directly help with appropriate configuration of certain parameters, such as the data distribution choice.

Second, small file and metadata accesses are the most frequent kind [13], [22]. The large files account for a very large fraction (80 percent) of storage space occupation while representing a very small percentage (10 to 20 percent) of the total number of files in a storage system [12]. In contrast, small files and metadata that are 4 KB or smaller account for the most user accesses [12], [23]. Thus, small files and file metadata that is very small in size should be stored with a replication-based scheme and large files should be stored with an erasure-code-based scheme for the performance and cost efficiency considerations. Knowing these two important workload characteristics is really important when designing the Cloud-of-Clouds storage system.

2.4 Replication versus Erasure Codes

Two common redundant data distribution methods used in Cloud-of-Clouds to achieve high availability of data are replication-based and erasure-code-based schemes. Although replication has the potential to increase availability and durability, it introduces two important challenges to system architects. First, system architects must increase the number of replicas to achieve high availability for large systems. For example, at least three replicas are required in Hadoop [24]. Second, the increased number of replicas introduces the extra bandwidth and storage overhead to the system,

TABLE 2
Comparison between Replication and Erasure Code Schemes

Schemes	Recovery	Performance	Cost
Replication	Easy	Low for large accesses	High
Erasure Codes	Hard	Low for small updates	Low

especially for large files. However, for the small files, replication is still an efficient way to provide the best performance with a small bandwidth and storage overhead. This is because small files only account for a tiny fraction of the bandwidth and storage capacity requirement, making replication on them profitable and productive considering the substantial performance benefits, both in the normal and recovery states.

An erasure code provides redundancy with much less space overhead than strict replication. Erasure codes divide an object into m fragments and recode them into a larger n fragments such that the original fragments can be recovered from a subset of the n fragments. The fraction $r = m/n$ is called the code rate. A rate r code increases the storage cost by a factor of $1/r$. For example, the RAID5 code can be described by an $(m=4, n=5)$ erasure code. The key property of erasure codes is that the original object can be reconstructed from any m fragments. The main advantage of erasure codes is the high space efficiency with good availability. However, since any m correctly verified fragments must be used to reconstruct a given lost fragment, it will introduce two serious performance problems. One is the extra time required to record the redundancy information, especially for small files. Take RAID5 for example, a small-file update will induce two read operations and two write operations. The other is the large amount of network traffic required to reconstruct data when a cloud provider suffers an outage or fails. For the RAID5 example, the recovery of a lost small file on the downed/failed storage provider will require read traffic from all surviving cloud storage providers. On the other hand, erasure codes offer a particular advantage for large files in that their access latency is reduced by virtue of the parallel accesses among multiple cloud providers.

Table 2 compares the two data distribution schemes, replication and erasure code. In general, replication provides better performance while erasure codes provide better storage efficiency. However, the former imposes extremely high bandwidth and storage overhead, while the latter does not provide the robustness and expected high access performance in the Cloud-of-Clouds particularly for large files. It therefore hints at the possibility of a certain combination of the two that tries to retain their respective advantages and while hiding their disadvantages to provide the most appropriate redundant data distribution scheme in Cloud-of-Clouds.

These important observations, combined with the urgent need to address the availability problem of cloud storage systems, motivate us to propose HyRD. In HyRD, large files are distributed in multiple cost-efficient cloud storage providers with erasure-coded data redundancy while small files and file system metadata are replicated on multiple high-performance cloud storage providers. By exploiting the workload

characteristics and the diversity of cloud providers, HyRD retains the desirable advantages of both the replication-based and erasure-code-based schemes to improve the availability of the Cloud-of-Clouds storage system.

3 THE DESIGN OF HYRD

In this section, we first outline the main design objectives of HyRD. Then we present its architecture overview, data distribution and cloud storage service evaluation methods, and the recovery scheme. Finally we present the data Consistency issue in HyRD

3.1 The Design Objectives of HyRD

The design of HyRD aims to achieve the following three objectives.

- *Improving the cloud storage dependability* - By redundantly distributing user data in a Cloud-of-Clouds, it solves the vendor lock-in problem. With data redundancy schemes of replication and erasure codes, the service unavailability problem caused by the outage of single individual cloud storage providers is avoided.
- *Reducing the user access latency* - By using replication for the small files and file system metadata, the performance issue of update operations is avoided. Moreover, by using erasure codes for the large files, the access latency is reduced by exploiting the access parallelism across multiple cloud storage providers.
- *Improving the cost efficiency* - Since HyRD uses erasure codes to store the large files that occupy the most storage capacity, the overall storage efficiency is improved. Moreover, while small files and file system metadata account for most user accesses, they occupy disproportionately small capacity. Thus, replication on them does not increase overall capacity cost noticeably.

3.2 HyRD Architecture Overview

Fig. 1 shows a system architecture overview of our proposed HyRD in the context of a Cloud-of-Clouds. Since more cloud storage services are provided by commercial cloud providers, the cloud providers are not allowed to execute users' code on the cloud storage side. As shown in Fig. 1, HyRD resides on the client side and interacts with the cloud storage providers via their standard interfaces without any modification. Thus, HyRD can be easily applied to any cloud storage providers to use their cloud storage services.

HyRD has three main functional modules: Workload Monitor, Request Dispatcher, and Cost & Performance Evaluator. The *Workload Monitor* module is responsible for classifying the incoming write data into file metadata, large files and small files. The qualification of a file being large or small is workload independent but related to the access latency. We have conducted performance evaluations to select the best threshold to determine a file's type in Section 4. Based on the data type information (i.e., file system metadata, small file, or large file), the *Request Dispatcher* module decides which redundancy scheme should be used for the incoming data, and distributes the data to the

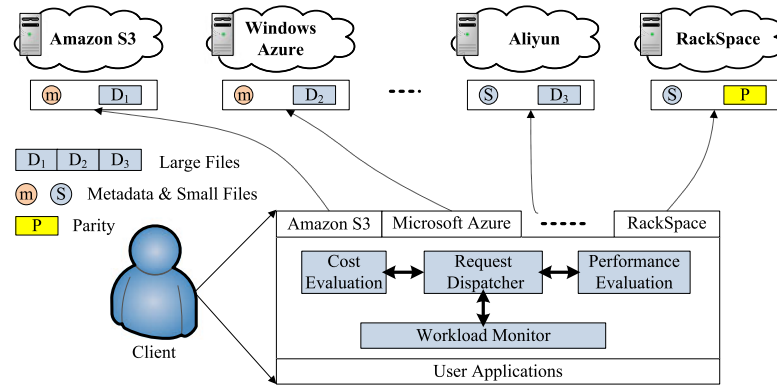


Fig. 1. System architecture of HyRD.

corresponding cloud storage providers. The *Cost & Performance Evaluator* module is responsible for evaluating the cloud storage services from the perspectives of cost and performance. The cost characteristics of the cloud storage providers are summarized in Table 1 in Section 4 and the performance characteristics are mainly described in terms of the access latency. These evaluation results will enable the Request Dispatcher module to select the appropriate cloud storage providers.

3.3 Data Distribution

Cloud-of-Clouds storage system is more reliable than their individual cloud storage provider. In order to build highly dependable and reliable systems out of less reliable parts, storage systems introduce redundancy. Currently, two redundancy schemes are widely involved: replication and erasure code. In replicated systems, objects are simply copied several times with each copy residing on a different physical device. While such an approach is simple and direct, more elaborate approaches such as erasure coding can achieve equivalent levels of data protection while using less redundancy. However, previous studies have shown that in a system with a large volume of active data, objects accessed frequently, and in which it is important to minimize latency, replication has the advantage. On the other hand, a system with mostly inactive data, archival objects accessed rarely, and primarily concerned with storage costs, would be better off using erasure coding.

Based on the observations in Section 2.3, the workload is mix of large and small I/O requests with different access patterns. Thus, HyRD uses hybrid data distribution scheme by exploiting the workload characteristics to choose either replication or erasure codes to distribute data among multiple cloud storage providers. At the present, HyRD only exploits the data type and file size characteristics. Besides file accesses, file system metadata blocks are critical to system performance. Before accessing a file, its metadata blocks must be loaded into the client memory. HyRD uses replication to store the file system metadata and groups the metadata in a directory together to exploit the access locality. For the file accesses, the file size is a critical parameter for HyRD to distribute data among multiple cloud storage providers. Since small files occupy a disproportionately small storage capacity and their updates are expensive in an erasure-code-based scheme, HyRD uses a replication-based scheme to store them. For large files that occupy a

disproportionally large storage capacity and need parallel accesses to improve their performance, HyRD uses an erasure-code-based scheme to store them. However, how to distinguish a large file from a small file is nontrivial as it sensitively depends on a file-size threshold. We have conducted sensitivity experiments to investigate the file-size threshold, as shown in Section 4.

The degree of data replication for file system metadata and small files within HyRD determines how resilient it is to cloud provider outages and failures, obviously the higher replication degree the more desirable. Unfortunately, higher degree of replication also comes at a higher cost both in terms of storage space and write/update latency. Thus, there is a trade-off among resiliency, cost and performance. For example, higher degree of replication (i.e., more replicas) imply higher resiliency but also lower performance for write/update operations for file system metadata and small files. A recent survey of the cloud service outages indicates that two concurrent cloud outages are extremely rare [25]. This, combined with the known fact that high degree of replication significantly degrades system performance while also incurring high space cost, makes it sensible to choose the replication level of 2 in our current HyRD design. Nevertheless, it must be noted that the degree of replication in HyRD is configurable to satisfy the requirements of different users.

3.4 Service Evaluation

The services offered by the cloud storage providers are diverse. Based on the design objectives of HyRD, it is important to evaluate the different cloud storage providers from the perspectives of performance and cost. In the HyRD evaluator module, the cloud storage providers are classified into two categories: performance-oriented providers where the data access latency is lower and cost-oriented providers where the storage capacity price is lower. For example, Aliyun is less expensive than Amazon S3 in storage price (per GB/month), as shown in Table 1. A particular cloud storage provider can be in one category or both, as shown in Fig. 2. Previous studies have shown that file system metadata and small files are accessed more frequently than large files [12]. Thus these data should be stored in performance-oriented providers for fast accesses. Because large files contribute to a disproportionately large storage capacity and thus the associated cost, HyRD puts them in the cost-oriented providers. To optimize performance of large files,

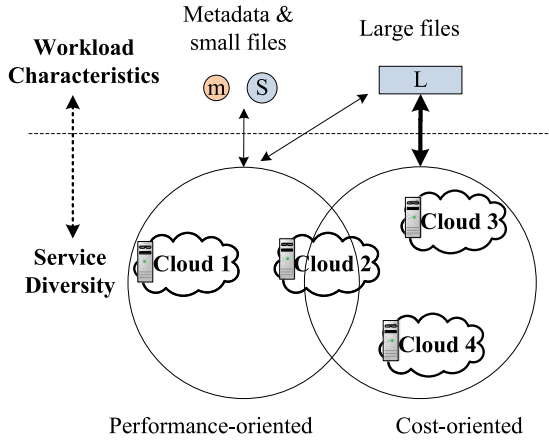


Fig. 2. Data distribution between the performance-oriented and cost-oriented providers.

some frequently accessed large files are also placed in performance-oriented providers, as shown in Fig. 2.

In our current implementation, the cost evaluation is configured manually by collecting the cost values from the official website of the cloud storage providers. The performance evaluation is performed online by periodically collecting the response times of different cloud storage providers. Because the cloud storage performance is affected by many factors, such as the network latency and the I/O intensity. To have a fair comparison and alleviate the influence of the other factors, HyRD extracts the performance values during system idle period for different cloud storage providers.

3.5 Recovery from Service Outage

An outage of cloud storage service is different from a disk failure in a disk array [26], [27], [28] in that the former results in a period of time during which cloud storage service is unavailable. The period may be hours and up to days which wreaked havoc on startups and other enterprises relying on cloud storage provider. For example, Microsoft Azure, which had a highly publicized cross-region outage in November of year 2014, had the most downtime at nearly 40 hours [5]. However, most outages will return to the normal state eventually. Thus, recovery in case of service outage in HyRD includes two phrases: (1) reconstruction on-demand during the unavailable period and (2) new written data update upon service's return to the normal state.

During the service unavailable period, all the write/update operations are performed as usual. For the update operations, the changes are logged; whereas, all the read operations are performed with on-demand read reconstruction, as shown in Fig. 3. For the file system metadata and small files, the requested data is directly fetched from the replicated providers. The large files are reconstructed using the erasure-code redundancy. The unrequested data on the unavailable provider is not reconstructed and migrated to other storage providers. During the service unavailable period, the data on the off-line cloud storage provider may be invalid due to the write/update operations. Upon the unavailable provider's return to the normal state, the recorded write/update logs will perform the consistency updates on the returned provider to make the data consistent. When the logs are completely processed, the recovery process completes.

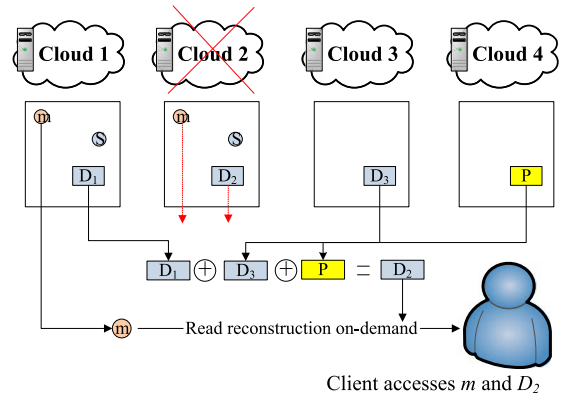


Fig. 3. The workflow of on-demand read reconstruction.

3.6 Data Consistency

Data consistency in our HyRD design includes the following two aspects: (1) The write data must be atomically stored on the cloud storage providers, (2) The user read requests must fetch the right data during the service unavailable period.

First, each write operation may involve multiple read/write operations in multiple cloud storage providers in HyRD. For example, a single small update will induce two read operations for the old data and parity blocks and two write operations for the new data and parity blocks in erasure-code-based scheme. Thus, completion of a single update will wait until all the induced read/write operations completed. To ensure the data consistency among the multiple cloud storage providers, HyRD will first buffer the write data. Upon all the induced operations are completed, the data in the buffer will be freed to store the new data blocks. Once the induced operations are not eventually performed and returned, the write operation will be performed again to make the data eventually written to the Cloud-of-Clouds storage system.

Second, read requested data may be on the unavailable cloud storage provider thus will cause a read failure. In HyRD, these read requests will be either redirected to the mirroring provider or reconstructed from all the other providers. Due to the data distribution scheme in HyRD, the read requests on the unavailable cloud storage provider will not incur significant extra I/O operations. For a small read request or metadata accesses, they will be serviced by the mirroring provider. For a large read request, since it already needs access all the cloud storage providers, the only induced extra I/O operation is the read request of the parity block which is needed to reconstruct the unavailable data block. Thus, the returned data of the user read requests to the unavailable cloud storage provider is guaranteed to be up-to-date.

4 PERFORMANCE EVALUATIONS

In this section, we first describe the prototype implementation and the experimental setup and methodology. Then we evaluate the performance of HyRD through extensive trace-driven experiments.

4.1 Prototype Implementation

Our prototype is built as an independent module on the client side. To interact with multiple cloud storage providers,

TABLE 3
Category Classifications of Different Cloud Storage Providers

Providers	Amazon S3 [2]	Windows Azure [15]	Aliyun [16]	RackSpace [17]
Category	Cost-oriented	Performance-oriented	Both	Cost-oriented

we have implemented a middleware of general cloud storage API, short for GCS-API. The GCS-API middleware hides the complexity of the cloud storage providers at the system level. Moreover, with such middleware, it is easy to add new cloud storage providers to the HyRD system.

Since each cloud storage service is modeled as a passive storage functional entity that supports five functions: List (lists the files of a container in the cloud), Get (reads a file), Create (creates a container), Put (writes or modifies a file in a container) and Remove (deletes a file). By passive storage functional entity, we mean that no operations other than what is needed to support the aforementioned five functions are executed. We assume that access control is provided by the system in order to ensure that read requests are only allowed to invoke the List and Get functions. To easily use the various cloud storage services, HyRD uses the REpresentational State Transfer APIs (short for RESTful APIs) to perform the operations. RESTful APIs are application program interfaces (APIs) that use HTTP requests to perform the above five functions. RESTful APIs explicitly take advantage of HTTP methodologies defined by the RFC 2616 protocol. Besides the above five functions, the Evaluation module in HyRD will directly interact with the individual cloud storage providers to evaluate the corresponding values.

4.2 Experimental Setup and Methodology

Our tests are conducted in a desktop PC (client) with an Intel i5-3470 3.2 GHz quad-core processor, with 4 GB of RAM and 1 Gigabit Ethernet connected to the China Education and Research Network [29]. Currently, our evaluations use the following four cloud storage providers in their default configurations: Amazon S3 [2], Windows Azure [15], Aliyun [16] and Rackspace [17]. The cost analysis is based on the values in Table 1 as of September 10th 2014. Table 3 shows the category classifications for the four major providers.

Because cost analysis is a long-term evaluation, similar to RACS, we used a trace-driven simulation to understand the costs associated with hosting large digital libraries in the cloud. Our trace covers one year of activity on the

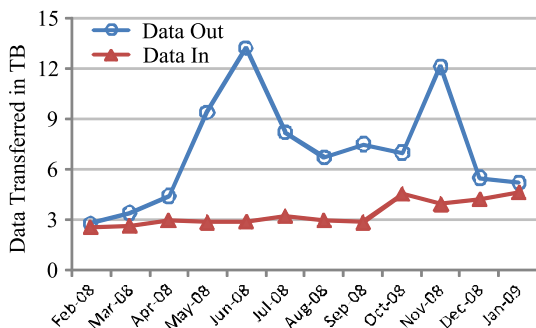
Internet Archive (IA) servers [30] from February 2008 to January 2009. Fig. 4 shows the amount of data written/read to/from the Internet Archive servers and the number of read/write requests issued to the Internet Archive servers during this one-year period. As shown in Fig. 4, the volume of data transferred is dominated by reads that outweigh writes by ratio of 2.1:1 and read requests outnumber write requests by a ratio of 3.5:1. The trace represents HTTP and FTP interactions that read and write various documents and media files (images, sounds, videos) stored at the Internet Archive and served to users. We believe that this trace is a good example of the type of workloads generated by online digital library systems, both in terms of the file sizes and request patterns.

To measure performance, we use the PostMark [31] benchmark tool to generate the file accesses as it is not practical to replay one year's trace. PostMark is designed to portray performance in desktop applications like electronic mail, netnews and web-based commerce, etc. We use PostMark to generate an initial pool of random text and image files ranging in size from a lower bound of 1,024 bytes to a higher bound of 100 M bytes. By default, we set the file-size threshold at 1 MB to distinguish large files from small files. For erasure-coded redundancy, we choose the RAID5 scheme in HyRD as a case study to fairly compare with the RACS approach.

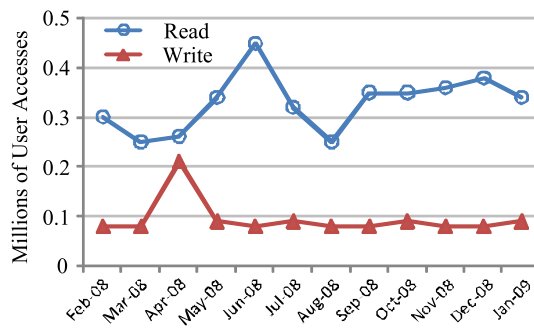
4.3 Cost Simulation and Analysis

In our cost simulation, it's assumed that the cloud services start with an empty storage without any data being pre-loaded. We estimate the cloud cost of moving the IA data to the cloud by using the up-to-date pricing schemes of the leading public cloud storage providers. Besides the bandwidth and storage costs, cloud providers also charge meta-data operations, such as Put, Copy, Post (short for 3Ps), List, Get and other operations based on per 10 K transactions, as shown in Table 1.

Fig. 5a shows the estimated monthly cost of servicing the Internet Archive by using single-cloud storage providers



(a) Data transferred in TB



(b) User read/write requests count

Fig. 4. The amount of data written/read to/from and the number of read/write requests issued to the Internet Archive servers as a function of time during a one-year period.

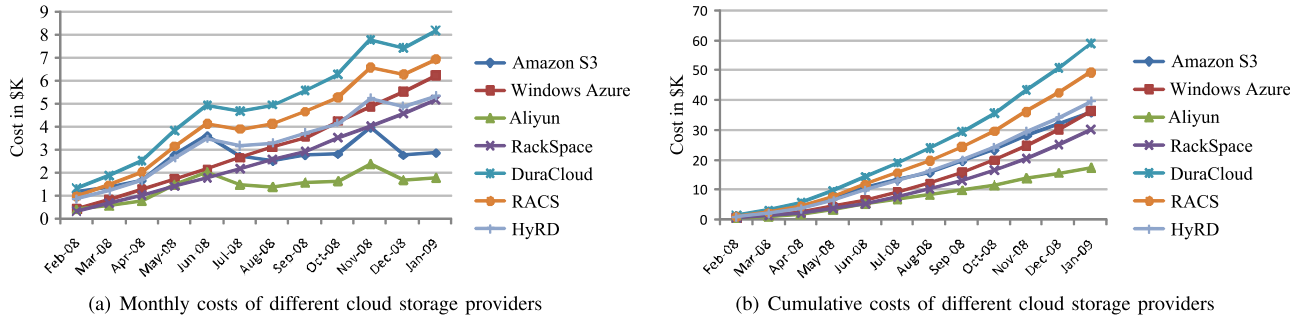


Fig. 5. Estimated monthly and cumulative costs of hosting storage services on the cloud for different schemes.

(i.e., Amazon S3, Windows Azure, Aliyun, and Rackspace), the DuraCloud scheme that fully replicates all data between two cloud storage providers, and the RACS scheme and our HyRD scheme that both distribute data redundantly among four cloud storage providers. While RACS uses the RAID5 scheme to distribute all data, HyRD relies on a hybrid replication and RAID5 scheme to dynamically and adaptively distribute data based on their type and size. From Fig. 5a, we can see that the monthly costs of all the schemes, except for Amazon S3 and Aliyun, increase nearly monotonously. The reason is that with each additional month, the monthly cost not only includes storage cost and read cost of the current month, but also includes the storage cost of all previously written data. However, for the Amazon S3 and Aliyun providers, while their storage costs are lower than Windows Azure and RackSpace by more than four times, their own read costs are much higher than their storage costs. This means that for Amazon S3 and Aliyun their monthly bills are dominated by the read costs that can fluctuate from month to month. In other words, the monthly costs for the Amazon S3 and Aliyun providers depend much more on the read (data-out) operations than on write (data-in) operations.

Fig. 5b shows the cumulative costs with different storage providers. First, we can see that DuraCloud is the most costly provider and Aliyun is the least costly provider. The high cost of DuraCloud comes from the full replication scheme that doubles the storage requirement and thus results in a storage cost that is the sum of those of the two involved individual providers. As expected, the cumulative storage cost increases from month to month as more data are accumulatively stored with time. Aliyun has the lowest cloud cost since it charges very little for the stored data and other operations, such as data out and metadata operations. Second, we see that the three Cloud-of-Clouds schemes (DuraCloud, RACS and HyRD) are more costly than the individual cloud storage providers. The reasons are twofold: (1) the Cloud-of-Clouds schemes add extra data redundancy that incurs additional storage cost; and (2) the update and write operations in the Cloud-of-Clouds will incur extra bandwidth cost due to the increased read operations. Third, the cloud cost of the HyRD scheme is 33.4 and 20.4 percent lower than that of the DuraCloud and RACS schemes, respectively. Compared with the DuraCloud scheme that uses full replication, both RACS and HyRD require less storage space overhead and thus achieve lower storage cost. Relative to the RACS scheme, our HyRD scheme requires less the storage cost by placing the many more large files in

the cost-oriented cloud storage providers, such as Aliyun and RackSpac. Moreover, by reading data from the cost-oriented cloud storage providers, HyRD's cloud cost due to the data out operations is also reduced. The breakdown of the total cloud cost into storage cost, read cost and other operations cost, shown in Fig. 6, further validates and illustrates the reasons behind HyRD's advantage over RACS.

Fig. 6 shows the breakdown of the cloud storage cost in the seven systems, i.e., four single-cloud storage providers and three Cloud-of-Clouds systems, where the DuraCloud system involves two single-cloud storage providers and the RACS and HyRD systems each involve four single-cloud storage providers. We can see that the cloud cost, or the user bill, is dominated by the storage cost for the Windows Azure, RackSpace, DuraCloud, RACS and HyRD schemes. It clearly shows that the costs charged for Windows Azure and RackSpace are purely storage cost. This is because that these two cloud storage providers only charge the storage cost, and make the other operations free for the users, as shown in Table 1. For Amazon S3 and Aliyun, on the contrary, it is the read traffic cost that dominates the user bill. The reason is obvious: the read operations in Amazon S3 and Aliyun are charged more than six and four times higher than the storage space, and the volume of data transferred out (read) is larger than that of data transferred in (written). On the other hand, even though DuraCloud chooses Aliyun as one of its replication providers, it is still the most expensive scheme. In contrast, by storing large files across multiple single-cloud providers, RACS and HyRD are able to reduce the storage cost significantly. In fact, the larger the number of single-cloud providers are involved, the more storage cost will be saved. Furthermore, by placing the large files in the cost-oriented providers, such as Aliyun, HyRD is able to cut more storage cost, making its storage cost much lower than that of RACS and only slightly higher than some

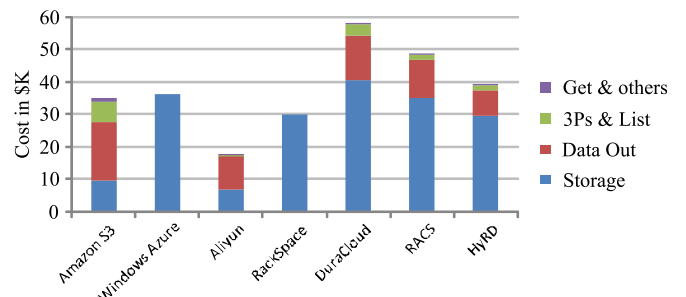


Fig. 6. Breakdown of cloud costs for different schemes.

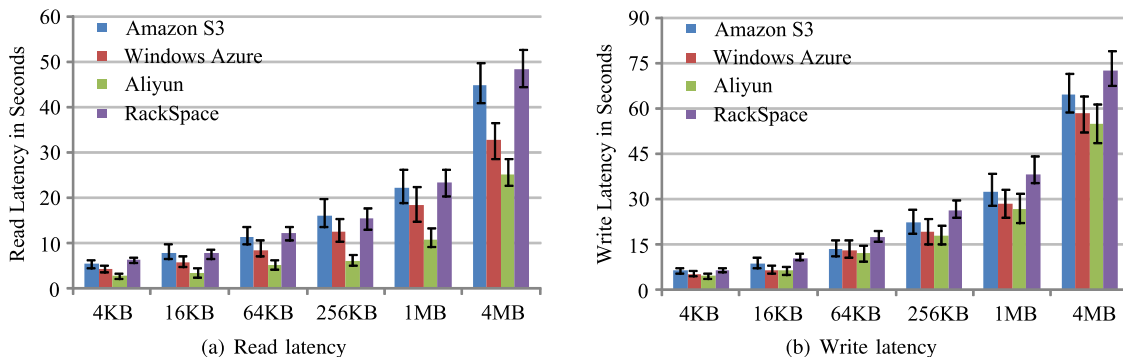


Fig. 7. Read/Write latency as a function of file size for single-cloud storage providers.

single-cloud providers, such as the Amazon S3 and Windows Azure schemes.

4.4 Performance Results

In order to understand the performance of HyRD in a real deployment, we use the PostMark benchmark tool to run several workloads accessing a Cloud-of-Clouds composed of four popular single-cloud providers of Amazon S3, Windows Azure, Aliyun and RackSpace. Since the Internet bandwidth is not stable from time to time, we run each experiment for three times and use the average latency results with the deviation values. These experiments took place during about three months between July 5, 2014 and October 7, 2014.

We first evaluate the performance of individual single-cloud storage providers as a function of the request sizes of 4 KB, 16 KB, 64 KB, 256 KB, 1 MB and 4 MB, as shown in Fig. 7. From the results presented in the figure, we can draw some interesting observations. First, Aliyun has the lowest access latency among the four single-cloud storage providers. This, combined with the fact that it has the lowest cloud cost as demonstrated in Figs. 5 and 6, makes Aliyun a unique cloud storage provider in that it is both performance-oriented and cost-oriented and thus explains its categorization in last row of Table 1. Second, there is a huge variance among the performance and the cost of the different cloud providers. It implies an important advantage of the Cloud-of-Clouds: we can exploit the workload characteristics and diversity of cloud storage providers to distribute data among multiple cloud storage providers. Third, when the file size increases from 1 to 4 MB, the access latency seems to increase disproportionately, which implies that the data transfer latency is disproportionately high at this level of file size and thus presents a clear gap the latency trend. Thus we set the file-size threshold at 1 MB to distinguish large files from small files.

Fig. 8 shows the benchmark results in terms of access latency of the different schemes both in the normal state and in the service outage period of any one of the single-cloud storage providers. In the evaluations, we configure PostMark to issue files of size ranging from 1 KB to 100 MB to simulate a desktop PC client. We set the performance of a single-cloud Amazon S3 storage provider as the baseline. In the normal state, we see that HyRD performs the best among all the schemes. Its access latency is 58.7 and 34.8 percent lower than the DuraCloud and RACS schemes, respectively. Both RACS and HyRD distribute large files across multiple

single-cloud storage providers, which enables them to exploit the access parallelism to improve the performance. However, for small files, the RACS scheme is less effective than HyRD. It is further validated by the performance results during the outage period when a single-cloud storage provider is off-line. For RACS, accessing (reading) the metadata or small files on the off-line provider will require it to access all the other three single-cloud storage providers to reconstruct the unavailable data. This will significantly increase the read traffic and decrease the bandwidth utilization. In contrast, the small-file/metadata access latency of neither DuraCloud or HyRD is noticeably affected by the service outage. The reason is that the metadata and small files are simply fetched from the surviving single-cloud storage provider that stores replicas of the unavailable data in the DuraCloud and HyRD schemes. In fact, the access latency of HyRD is 46.3 percent lower than that of RACS during the service outage period. Moreover, upon a service outage, the access latency of DuraCloud is better than that in the normal state since no double writes or updates are performed. This explains why the access latency of HyRD is only 27.3 percent lower than that of DuraCloud during the service outage period.

4.5 Estimates of Cost-Effectiveness

To reasonably estimate and quantify the cost-effectiveness of HyRD in comparison to the state-of-the-art Cloud-of-Clouds schemes and single-cloud schemes, we use the cost-latency product as a measure for cost-effectiveness, taking inspiration from the energy-latency product that is commonly used in the computer architecture literature to quantify energy

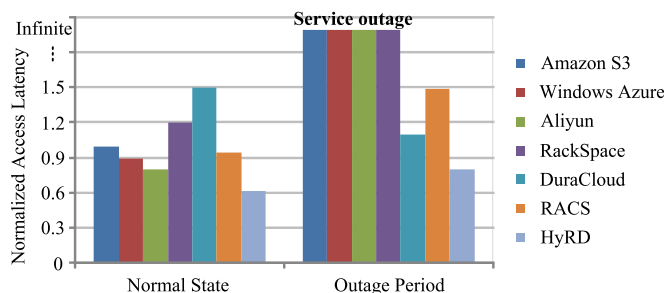


Fig. 8. Access latency results of benchmark-driven experiments on real cloud deployment of the different schemes both in the normal state and in the service outage period of any one of the single-cloud storage providers. Note that the results are normalized to Amazon S3 and we set the Windows Azure service off-line to emulate its outage when evaluating the three Cloud-of-Clouds schemes.

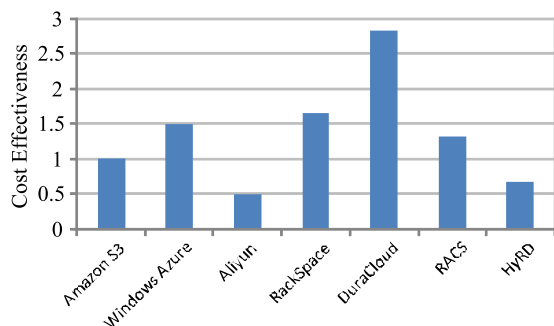


Fig. 9. Benchmark results on cost effectiveness of the different schemes. Note that the lower the value, the better for the cost effectiveness.

efficiency [32]. The lower the cost-latency product value of a scheme, the more cost-effective the scheme is.

Fig. 9 shows cost-effectiveness, in terms of the cost-latency product, of the different schemes based on the latency and cost results obtained from the benchmark-driven experiments presented earlier in this section. We can see that among all the schemes, Aliyun is the most cost effective, followed by HyRD. As shown in Table 1, Aliyun is therefore categorized into both a performance-oriented and cost-oriented cloud storage provider. However, as we explained in Section 2.1, all single-cloud storage providers, including Aliyun, have the inherent vendor lock-in problem that Cloud-of-Clouds is designed to address. It is for this reason we believe that it is more meaningful to compare among the Cloud-of-Clouds schemes for their cost-effectiveness. Among the three Cloud-of-Clouds schemes under study HyRD is clearly the most cost effective, outperforming DuraCloud and RACS by 76.2 and 49.4 percent respectively, and achieves comparable or better cost-effectiveness than that of the single-cloud storage providers except for Aliyun. The reasons behind HyRD's superiority in cost-effectiveness are two-fold. First, by placing large files in cost-oriented providers, the storage cost is reduced. Second, by placing frequently accessed file system metadata and small files in performance-oriented providers, the overall access latency is reduced.

5 RELATED WORK

As cloud storage becomes popular and cost efficient, more and more organizations and individual users will move their data to the cloud. Besides performance and security, availability of the cloud storage service is becoming increasingly more important for users. The notion of Cloud-of-Clouds is an effective approach to addressing the availability issue caused by the service outages of single-cloud storage providers.

There are several systems proposed for Cloud-of-Clouds [3], [4], [7], [8], [33], [34], [35], [36], [37], [38]. RACS [3]

uses erasure coding to mitigate the vendor lock-in problem encountered by a user when switching cloud vendors. It transparently stripes data across multiple cloud storage providers with RAID-like techniques used by disks and file systems. HAIL [33] provides integrity and availability guarantees for stored data. It allows a set of servers to prove to a client that a stored file is intact and retrievable by the approaches adopted from the cryptographic and distributed-systems communities. NCCloud [7] achieves cost-effective repair for a permanent single-cloud provider failure to improve availability of cloud storage services. It is built on top of network-coding-based storage schemes called regenerating codes with an emphasis on storage repair, excluding the failed cloud in repair.

The above three systems are all based on erasure codes or network codes. In contrast, DuraCloud [14] utilizes replication to copy user content onto several different cloud storage providers to provide better availability. Moreover, it ensures that all copies of user content remain synchronized. However, users will pay more money for the additional storage space and bandwidth required by DuraCloud. DEPSKY [4] improves the availability and confidentiality of commercial storage cloud services by building a Cloud-of-Clouds on top of a set of storage clouds, combining Byzantine quorum system protocols, cryptographic secret sharing, replication and the diversity provided by the use of several cloud providers. Besides these optimizations, there are also some studies targeted for data security improvement across multiple cloud storage providers, such as Hybris [35] and CDStore [34]. Though they have different objectives, they are orthogonal to and can be incorporated with the existing Cloud-of-Clouds schemes. Table 4 summarizes the state-of-the-art redundant data distribution schemes for availability improvement. Different from these approaches, our proposed HyRD scheme takes the workload characteristics and the diversity of cloud storage providers, specially the file sizes, into the design of the redundant data distribution strategy so that the advantages of both the replication and erasure codes are exploited while hiding their disadvantages. As a result, both performance and storage efficiency are improved with the availability guarantee.

Integrating replication and erasure codes into one system is not a new idea. Our proposed HyRD takes inspirations from previous studies in the data organizations for RAID and file systems [39], [40], [41], [42]. For different RAID levels [43], replication-based disk array (RAID1) and parity-based disk arrays (RAID4/5) provide different performance and storage efficiency. HP AutoRAID [39] provides a two-level storage hierarchy inside a monolithic disk array controller. In the upper level of this hierarchy, RAID1 provides full redundancy and better performance. In the lower level,

TABLE 4
Comparison Between HyRD and the State-of-the-Art Schemes

Schemes	Redundancy	Recovery	Performance	Cost
RACS [3]	Erasure Codes	Hard	Low for small updates	Low
DuraCloud [1], [14]	Replication	Easy	Low for large accesses	High
DepSky [4]	Replication	Easy	Low for large accesses	High
NCCloud [7]	Network Codes	Moderate	Low for small updates	Low
HyRD	Replication and erasure code	Easy	High	Low

RAID5 parity protection is used to achieve lower storage cost. It automatically and transparently manages migration of data blocks between these two levels as access patterns change. Hot Mirroring [40] similarly combines RAID1 and RAID5 layouts, keeping hot data in the RAID1 portion and cold data in the RAID5 portion. It is a single-box solution and uses metadata to control the placement of data among disks comprising the disk array. In contrast to HP AutoRAID and Hot Mirroring, HyRD exploits the workload characteristics and the heterogeneity of cloud providers to choose between the replication redundancy and the erasure-coded redundancy to distribute data among multiple cloud storage providers, thus improving the availability of cloud storage services from the user's perspective.

6 CONCLUSION

Availability of cloud storage services is one of the main factors that the users must consider seriously when deciding whether or not to move their data to the cloud. Depending on a single cloud storage provider has the inherent vendor lock-in problem that can potentially cost the user dearly. This paper proposed a hybrid redundant data distribution approach, called HyRD, by exploiting the workload characteristics and the diversity of cloud storage providers to improve the storage availability in Cloud-of-Clouds. In HyRD, large files are distributed in multiple cost-oriented cloud storage providers with the erasure-coded data redundancy while small files and file system metadata are replicated on multiple performance-oriented cloud storage providers. By exploiting the workload characteristics and the heterogeneity of cloud providers, both the advantages of erasure codes and replication are exploited while their disadvantages are alleviated. The experiments conducted on our lightweight prototype implementation of HyRD show that HyRD significantly outperforms existing Cloud-of-Clouds schemes, such as RACS and DuraCloud, in terms of the I/O performance and cost-effectiveness.

HyRD is an ongoing research project and we are currently exploring several directions for future research. First, we will apply data deduplication in the HyRD module to eliminate the redundant data and reduce the total data transferred over the network, thus further improving the performance and cost efficiency [44], [45]. However, data deduplication requires powerful computing resources and extra memory space while HyRD is located in the client side. Applying data deduplication in HyRD is not easy and needs careful design considerations. Second, we will extend the HyRD design to consider the specific features of the diverse cloud storage services, thus further improving the flexibility of HyRD and the efficiency of cloud storage services.

ACKNOWLEDGMENTS

This work is supported by the National Natural Science Foundation of China under Grant No. 61100033, No. 61472336 and No. 61402385, the US National Science Foundation (NSF) under Grant No. NSF-CNS-1116606 and NSF-CNS-1016609, National Key Technology R&D Program Foundation of China (2015BAH16F02) and Fundamental Research Funds for the Central Universities (No. 20720140515). S. Wu is the corresponding author.

REFERENCES

- [1] E. Allen and C. M. Morris. (2009, Jul.). Library of congress and duracloud launch pilot program using cloud technologies to test perpetual access to digital content, in *Proc. Library Congress, News Release* [Online]. Available: <http://www.loc.gov/today/pr/2009/09-140.html>
- [2] (2014). Amazon S3 [Online]. Available: <http://aws.amazon.com/s3/>
- [3] H. Abu-Libdeh, L. Princehouse, and H. Weatherspoon, "RACS: A case for cloud storage diversity," in *Proc. ACM Symp. Cloud Comput.*, Jun. 2010, pp. 229–240.
- [4] A. Bessani, M. Correia, B. Quaresma, F. André, and P. Sousa, "DepSky: Dependable and secure storage in a cloud-of-clouds," in *Proc. 6th Eur. Conf. Comput. Syst.*, Apr. 2011, pp. 31–46.
- [5] (2014). The Worst Cloud Outages of 2014 [Online]. Available: <http://www.cio.com/article/2597775/cloud-computing/162288-The-worst-cloud-outages-of-2014-so-far.html>
- [6] (2014). Tim Cook Says Apple to Add Security Alerts for iCloud Users, the Wall Street Journal [Online]. Available: <http://online.wsj.com/articles/tim-cook-says-apple-to-add-security-alerts-for-icloud-users-1409880977>
- [7] Y. Hu, H. Chen, P. Lee, and Y. Tang, "NCCloud: Applying network coding for the storage repair in a cloud-of-clouds," in *Proc/10th USENIX Conf. File Storage Technol.*, Feb. 2012, pp. 265–272.
- [8] B. Mao, H. Jiang, and S. Wu, "Improving storage availability in cloud-of-clouds with hybrid redundant data distribution," in *Proc. 29th IEEE Int. Parallel Distrib. Process. Symp.*, May 2015, pp. 633–642.
- [9] K. Rashmi, N. Shah, D. Gu, H. Kuang, D. Borthakur, and K. Ramchandran, "A "Hitchhiker's" guide to fast and efficient data reconstruction in erasure-coded data centers," in *Proc. ACM SIGCOMM Conf. Appl., Technol., Archit. Protocols Comput. Commun.*, Aug. 2014, pp. 331–342.
- [10] K. Rashmi, N. Shah, D. Gu, H. Kuang, D. Borthakur, and K. Ramchandran, "A solution to the network challenges of data recovery in erasure-coded distributed storage systems: A study on the facebook warehouse cluster," in *Proc. 5th USENIX Workshop Hot Topics File Storage Technol.*, Jun. 2013, pp. 1–5.
- [11] R. Villars, C. Olofson, and M. Eastwood, "Big data: What it is and why you should care, white paper, IDC," Jun. 2011, http://www.admin-magazine.com/HPC/content/download/5604/49345/file/IDC_Big_Data_whitepaper_final.pdf.
- [12] N. Agrawal, W. J. Bolosky, J. R. Douceur, and J. R. Lorch, "A five-year study of file-system metadata," in *Proc. 5th USENIX Conf. File Storage Technol.*, Feb. 2007, pp. 31–45.
- [13] A. Traeger, E. Zadok, N. Joukov, and C. Wright, "A nine year study of file system and storage benchmarking," *ACM Trans. Storage*, vol. 48, no. 2, pp. 1–56, 2008.
- [14] (2014). DuraCloud Project [Online]. Available: <http://www.duracloud.org/>
- [15] (2014). Windows Azure [Online]. Available: <http://www.windowsazure.cn/zh-cn/>
- [16] (2014). Aliyun Open Storage Service [Online]. Available: <http://www.aliyun.com/>
- [17] (2014). RackSpace [Online]. Available: <http://www.rackspace.com/cn/>
- [18] J. Gahm and J. Mcknight, "Medium-size business server & storage priorities," *Enterprise Strategy Group*, Jun. 2008, <http://www.esg-global.com/research-reports/medium-size-business-server-storage-priorities/>.
- [19] S. Wu, H. Jiang, D. Feng, L. Tian, and B. Mao, "Improving availability of RAID-structured storage systems by workload outsourcing," *IEEE Trans. Comput.*, vol. 60, no. 1, pp. 64–79, Jan. 2011.
- [20] (2014). EMC Presents Disaster Recovery Survey 2013 [Online]. Available: <http://emea.emc.com/microsites/2011/emc-brs-survey/index.htm>
- [21] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. H. Katz, A. Konwinski, G. Lee, D. A. Patterson, A. Rabkin, and M. Zaharia, "Above the clouds: A berkeley view of cloud computing," EECS Dept., Univ. of California, Berkeley, CA, USA, Tech. Rep. No. USB/EECS-2009-28, 2009.
- [22] Y. Wang, L. Alvisi, and M. Dahlin, "Gnothi: Separating data and metadata for efficient and available storage replication," in *Proc. USENIX Annu. Tech. Conf.*, Jun. 2012, pp. 413–424.

- [23] J. Lofstead, M. Polte, G. Gibson, S. Klasky, K. Schwan, R. Oldfield, M. Wolf, and Q. Liu, "Six degrees of scientific data: Reading patterns for extreme scale science IO," in *Proc. 20th Int. Symp. High Perform. Distrib. Comput.*, Jun. 2011, pp. 49–60.
- [24] Hadoop [Online]. Available: <http://hadoop.apache.org/>, 2014.
- [25] O. Khan, R. Burns, J. S. Plank, W. Pierce, and C. Huang, "Rethinking erasure codes for cloud file systems: Minimizing I/O for recovery and degraded reads," in *Proc. 10th USENIX Conf. File Storage Technol.*, Jan. 2012, pp. 251–264.
- [26] S. Wu, H. Jiang, D. Feng, L. Tian, and B. Mao, "WorkOut: I/O workload outsourcing for boosting the RAID reconstruction performance," in *Proc. 7th USENIX Conf. File Storage Technol.*, Feb. 2009, pp. 239–252.
- [27] S. Wu, H. Jiang, and B. Mao, "IDO: Intelligent data outsourcing with improved RAID reconstruction performance in large-scale data centers," in *Proc. 26th USENIX Large Installation Syst. Admin.*, Dec. 2012, pp. 17–32.
- [28] S. Wu, H. Jiang, and B. Mao, "Proactive data migration for improved storage availability in large-scale data centers," *IEEE Trans. Comput.*, vol. 64, no. 9, pp. 2637–2651, Sep. 2015.
- [29] (2014). China Education and Research Network [Online]. Available: http://www.edu.cn/english_1369/index.shtml
- [30] Internet Archive [Online]. Available: <http://www.archive.org/>, 2009.
- [31] (2010). Filesystem Benchmarking with PostMark from NetApp [Online]. Available: <http://www.shub-internet.org/brad/FreeBSD/postmark.html>
- [32] A. Banerjee, P. Wolkotte, R. Mullins, S. Moore, and G. Smit, "An energy and performance exploration of network-on-chip architectures," *IEEE Trans. Very Large Scale Integration Syst.*, vol. 17, no. 3, pp. 319–329, Mar. 2009.
- [33] K. D. Bowers, A. Juels, and A. Oprea, "HAIL: A high-availability and integrity layer for cloud storage," in *Proc. 16th ACM Conf. Comput. Commun. Security*, Nov. 2009, pp. 187–198.
- [34] M. Li, C. Qin, and P. P. C. Lee, "CDStore: Toward reliable, secure, and cost-efficient cloud storage via convergent dispersal," in *Proc. USENIX Annu. Tech. Conf.*, Jul. 2015, pp. 111–124.
- [35] D. Dobre, P. Viotti, and M. Vukolić, "Hybris: Robust hybrid cloud storage," in *Proc. ACM Symp. Cloud Comput.*, Nov. 2014, pp. 1–14.
- [36] D. Li, X. Liao, H. Jin, B. Zhou, and Q. Zhang, "A new disk I/O model of virtualized cloud environment," *IEEE Trans. Parallel Distrib. Syst.*, vol. 24, no. 6, pp. 1129–1138, Jun. 2013.
- [37] X. Zhang, M. Tsugawa, Y. Zhang, H. Song, C. Cao, G. Huang, and J. Fortes, "Towards model-defined cloud of clouds," in *Proc. 17th Int. Conf. Model Driven Eng. Lang. Syst.*, Sep. 2014, pp. 41–45.
- [38] R. Zhou, H. Chen, and T. Li, "Towards lightweight and swift storage resource management in big data cloud era," in *Proc. 29th ACM Int. Conf. Supercomput.*, Jun. 2015, pp. 133–142.
- [39] J. Wilkes, R. Golding, C. Staelin, and T. Sullivan, "The HP Auto-RAID hierarchical storage system," *Operating Syst. Rev.*, vol. 29, no. 5, pp. 96–108, 1995.
- [40] K. Mogi and M. Kitsuregawa, "Hot mirroring: A method of hiding parity update penalty and degradation during rebuilds for RAID5," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, Jun. 1996, pp. 183–194.
- [41] B. Mao, H. Jiang, S. Wu, L. Tian, D. Feng, J. Chen, and L. Zeng, "HPDA: A hybrid parity-based disk array for enhanced performance and reliability," *ACM Trans. Storage*, vol. 8, no. 1, pp. 1–20, 2012.
- [42] Y. Ma, T. Nandagopal, K. Puttaswamy, and S. Banerjee, "An ensemble of replication and erasure codes for cloud file systems," in *Proc. 32nd IEEE Int. Conf. Comput. Commun.*, Apr. 2013, pp. 1276–1284.
- [43] D. Patterson, G. Gibson, and R. Katz, "A case for redundant arrays of inexpensive disks (RAID)," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, Jun. 1988, pp. 109–116.
- [44] B. Mao, H. Jiang, S. Wu, Y. Fu, and L. Tian, "Read performance optimization for deduplication-based storage systems in the cloud," *ACM Trans. Storage*, vol. 10, no. 2, pp. 1–22, 2014.
- [45] B. Mao, H. Jiang, S. Wu, and L. Tian, "POD: Performance oriented I/O deduplication for primary storage systems in the cloud," in *Proc. 28th IEEE Int. Parallel Distrib. Process. Symp.*, May 2014, pp. 767–776.



Bo Mao received the BE degree in computer science and technology in 2005 from Northeast University; and the PhD degree in computer architecture in 2010 from the Huazhong University of Science and Technology. His research interests include storage system, Cloud computing, and Big Data. He is an assistant professor at the Software School of Xiamen University. He has more than 30 publications in international journals and conferences including *IEEE Transactions on Computers*, *IEEE Transactions on Parallel and Distributed Systems*, *ACM Transactions on Storage*, USENIX FAST, IPDPS, Cluster, USENIX LISA, MASCOTS, and ICPADS. He is a member of the IEEE, ACM, and USENIX.



Suzhen Wu received the BE and PhD degrees in computer science and technology and computer architecture from the Huazhong University of Science and Technology, in 2005 and 2010, respectively. She is an associate professor at the Computer Science Department of Xiamen University since August 2014. Her research interests include computer architecture and storage system. She has more than 30 publications in journal and international conferences including *IEEE Transactions on Computers*, *IEEE Transactions on Parallel and Distributed Systems*, *ACM Transactions on Storage*, USENIX FAST, USENIX LISA, IPDPS, MASCOTS, and ICPADS. She is a member of the IEEE ACM.



Hong Jiang received the BSc degree in computer engineering in 1982 from the Huazhong University of Science and Technology, Wuhan, China, the MASc degree in computer engineering in 1987 from the University of Toronto, Toronto, Canada, and the PhD degree in computer science in 1991 from the Texas A&M University, College Station, Texas, TX. He is currently a chair and Nedderman professor of the Computer Science and Engineering Department at the University of Texas at Arlington. Prior to joining UTA, he served as a program director at National Science Foundation (2013.1-2015.8) and he was at the University of Nebraska-Lincoln since 1991, where he was Willa Cather professor of computer science and engineering. He has graduated 13 PhD students who upon their graduations either landed academic tenure-track positions in PhD-granting US institutions or were employed by major US IT corporations. His present research interests include computer architecture, computer storage systems and parallel I/O, high-performance computing, big data computing, cloud computing, and performance evaluation. He recently served as an associate editor of the *IEEE Transactions on Parallel and Distributed Systems*. He has more than 200 publications in major journals and international Conferences in these areas, including *IEEE Transactions on Parallel and Distributed Systems*, *IEEE Transactions on Computers*, *ACM Transactions on Architecture and Code Optimization*, *Journal of Parallel and Distributed Computing*, ISCA, MICRO, USENIX ATC, FAST, LISA, ICDCS, IPDPS, Middleware, OOPLAS, ECOOP, SC, ICS, HPDC, INFOCOM, ICPP, etc., and his research has been supported by US National Science Foundation (NSF), DOD, and the State of Nebraska. He is a fellow of the IEEE and a member of ACM.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.