

# Data Storage Management in Cloud Environments: Taxonomy, Survey, and Future Directions

Yaser Mansouri, Adel Nadjaran Toosi,  
and Rajkumar Buyya, The University of Melbourne, Australia.

Storage as a Service (SaaS) is a vital component of cloud computing by offering the vision of a virtually infinite pool of storage resources. It supports a variety of cloud-based data store classes in terms of availability, scalability, ACID (Atomicity, Consistency, Isolation, Durability) properties, data models, and price options. Application providers deploy these storage classes across different cloud-based data stores not only to tackle the challenges arising from reliance on a single cloud-based data store but also to obtain higher availability, lower response time, and more cost efficiency. Hence, in this paper, we first discuss the key advantages and challenges of data-intensive applications deployed within and across cloud-based data stores. Then, we provide a comprehensive taxonomy that covers key aspects of cloud-based data store: data model, data dispersion, data consistency, data transaction service, and data management cost. Finally, we map various cloud-based data stores projects to our proposed taxonomy to validate the taxonomy and identify areas for future research.

Categories and Subject Descriptors: C.4 [Performance of Systems]: Reliability, availability, and serviceability; C.2.4 [Computer-Communication Networks]: Distributed Systems

General Terms: Data Management, Data Storage, and Data Replication

Additional Key Words and Phrases: Data Management, Data Storage, Data Replication, Data Consistency, Transaction Service, and Data Management Cost

## ACM Reference Format:

Yaser Mansouri, Adel Nadjaran Toosi, and Rajkumar Buyya, 2015. Data Storage Management in Cloud Environments: Taxonomy, Survey, and Future Directions. *ACM Comput. Surv.* 0, 0, Article 01 (May 2015), 37 pages.

DOI: <http://dx.doi.org/10.1145/0000000.0000000>

## 1. INTRODUCTION

The explosive growth of data traffic driven by social networks, e-commerce, enterprises, and other data sources has become an important and challenging issue for IT enterprises. This growing speed is doubling every two years and augments 10-fold between 2013 and 2020- from 4.4 ZB to 44 ZB. The challenges posed by this growth of data can be overcome with aid of using cloud computing services. Cloud computing offers the illusion of infinite pool of highly reliable, scalable, and flexible computing, storage, and network resources in a pay-per-use manner. These resources are typically categorized as Infrastructure as a Service (IaaS), where Storage as a Service (SaaS) forms one of its critical components.

SaaS provides a range of cloud-based data stores (*data stores* for short) that differs in data model, data consistency semantic, data transaction support, and price model. A popular class of data stores, called Not only SQL (NoSQL), has emerged to host applications that require high scalability and availability without having to support the ACID properties of relational database (RDB) systems. This class of data stores – such as PNUTS [Cooper et al. 2008] and Dynamo [DeCandia et al. 2007] – typically partitions data to provide scalability and replicates the partitioned data to achieve high availability. *Relational data store*, as another class of data stores, provides full-fledged relational data model to support ACID properties, while it is not as scalable as NoSQL data store. To strike a balance between these two classes, NewSQL data store was introduced. It captures the advantages of both NoSQL and relational data stores and initially was exploited in Spanner [Corbett et al. 2013].

To take the benefits of these classes, application providers store their data either in a single or multiple data stores. A single data store offers the proper availability, durability, and scalability. But reliance on a single data store has risks like vendor lock-in, economic failure (e.g., a surge in price), and unavailability as outages occur, and probably leads to data loss when an environmental catastrophe happens [Borthakur et al. 2011]. Geo-replicated data stores, on the other hand, mitigate these risks and also provide several key benefits. First, the application providers can serve users from the best data store to provide adequate

---

Author's addresses: Yaser Mansouri, Adel Nadjaran Toosi, and Rajkumar Buyya, are with the Cloud Computing and Distributed Systems (CLOUDS) Laboratory, School of Computing and Information Systems, The University of Melbourne, Australia.  
Email: [yase@student.unimelb.edu.au](mailto:yase@student.unimelb.edu.au), [{anadjaran,rbuyya}@unimelb.edu.au](mailto:{anadjaran,rbuyya}@unimelb.edu.au)

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2015 ACM. 0360-0300/2015/05-ART01 \$15.00

DOI: <http://dx.doi.org/10.1145/0000000.0000000>

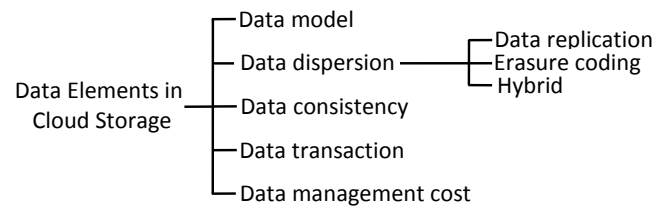


Fig. 1: Data elements in cloud storage

responsiveness since data is available across data stores. Second, the application can distribute requests to different data stores to achieve load balance. Third, data recovery can be possible when natural disaster and human-induced activities happen. However, the deployment of a single or multiple data stores causes several challenges depending on the characteristics of data-intensive applications.

Data-intensive applications are potential candidates for deployment in the cloud. They are categorized into transactional (ref. as *online transaction processing* (OLTP)) and analytical (ref. as *online analytical processing* (OLAP)) that demand different requirements. OLTP applications embrace different consistency semantics and are adaptable with row-oriented data model, while OLAP applications require rich query capabilities and compliance with column-oriented data model. These requirements are faced with several challenges, and they mandate that we investigate the key elements of data management in data stores as shown in Fig.1. The first five key elements are mature topics in the context of distributed systems and require how to apply them to data stores with possible modifications and adoptions if needed. The last element, *data management cost*, is a new feature for cloud storage services and it is important for users to optimize it while their Service Level Agreements (SLAs) are guaranteed.

The first element is *data model* that reflects how data is stored in and retrieved from data stores. The second element is *data dispersion* with three schemes. *Data replication* scheme improves availability and locality by moving data close to the user, but it is costly due to usually storing three replicas for each object in data stores. *Erasure coding* scheme alleviates this overhead, but it requires structural design to reduce the time and cost of recovery. To make a balance between these benefits and shortcomings, a combination of both schemes is exploited. We clarify these schemes and identify how they influence availability, durability, and user-perceived latency.

Other elements of data management are *data consistency* and *data transaction* that refer to coordination level between replicas within and across data stores. Based on CAP theorem [Gilbert and Lynch 2002], it is impossible to jointly attain Consistency, Availability, and Partition tolerance (referred to the failure of a network device) in distributed systems. Thus, initially data stores provide *eventual consistency*— all replicas eventually converge to the last updated value— to achieve two of three these properties: availability and partition tolerance. Eventual consistency is sometime acceptable, but not for some applications (e.g., e-commerce) that demand strong consistency in which all replicas receive requests in the same order. To obtain strong consistency, the recent endeavours bring transactional isolation levels in NoSQL/NewSQL data stores at the cost of extra resources and higher response time.

The last element is *cost optimization of data storage management* as the key driver behind the migration of application providers into the cloud that offers a variety of storage and network resources with different prices. Thus, application providers have many opportunities for *cost optimization* and *cost trade-offs* such as storage vs. bandwidth, storage vs. computing, and so on.

The main contributions of this paper are as follows:

- Comparison between cloud-based data stores and related data-intensive networks,
- Delineation on goals and challenges of intra- and inter-cloud storage services, and determination of the main solutions for each challenges,
- Discussion on data model taxonomy in three aspects: *data structure*, *data abstraction*, and *data access model*; and comparison between different levels/models of each aspect,
- Providing a taxonomy for different schemes of data dispersion, and determining when (according to the specifications of workload and the diversity of data stores) and which scheme should be used,
- Elaboration on different levels of consistency and determination on how and which consistency level is supported by the state-of-the-art projects,
- Providing a taxonomy for transactional data stores, classifying them based on the provided taxonomy, and analysing their performance in terms of throughput and bottleneck at message,
- Discussion on the cost optimization of storage management, delineation on the potential cost trade-offs in data stores, and classifying the existing projects to specify the research venue for future.

This survey is divided into seven sections. Section 2 compares cloud-based data stores to other distributed data-intensive networks and then discusses the architecture, goals and challenges of a single and multiple cloud-based data stores deployments. Section 3 describes a taxonomic of data model and Section 4 discusses different schemes of *data replication*. Section 5 elaborates on data consistency in

## Relational Schema:

Customer (CustomerId, ...),  
 Campaign (CampaignId, CustomerId, ...)  
 AdGroup(AdGroupId, CampaignId, ...)

## NoSQL Schema (key-value)

Customer (Key, Value<sub>1</sub>, Value<sub>2</sub>, ..., Value<sub>k</sub>)  
 Sample:

Customer (Cid1, CName1, CDegree1, CAddress1)  
 Customer (Cid2, CName2, CAddress2)  
 Customer (Cid3, CName3, CAddress3, CBalance3)

Fig. 2: Relational and NoSQL schemas

## NewSQL Schema:

Customer (CustomerId, ...),  
 Campaign (CampaignId, CustomerId, ...)  
 AdGroup(CustomerId, CampaignId, AdGroupId, ...)

Sample

Directory 1

Directory 2

Customer(1,...)  
 Campaign(1,3,...)  
 AdGroup(1,3,6,...)  
 AdGroup(1,3,7,...)  
 Campaign(1,4,...)  
 AdGroup(1,4,8,...)

Customer(2,...)  
 Campaign(2,5,...)  
 AdGroup(2,5,9....)

Cluster

Fig. 3: NewSQL Schema [Shute et al. 2013]

terms of *level*, *metric*, and *model*. Section 6 details the taxonomy of transactional data stores and Section 7 presents the cost optimization of storage management. Section 8 finally concludes the paper and presents some open research questions in the context of data storage management.

## 2. OVERVIEW

This section discusses a comparison between cloud-based data stores and other data-intensive networks (§2.1), the terms used throughout this survey (§2.2), the characteristics of data-intensive applications deployed in data stores (§2.3), and the main goals and challenges of a single and multiple data stores leveraged to manage these applications (§2.4).

### 2.1. A Comparison of Data-intensive Networks

Table I highlights similarities and differences in characteristics and objectives between cloud-based data stores and (i) Data Grid in which storage resources are shared between several industrial/educational organizations as Virtual Organization (VO), (ii) Content Delivery Network (CDN) in which a group of servers/datacenters are located in several geographical locations to serve users contents (i.e. application, web, or video) faster, and (iii) Peer-to Peer (P2P) in which a peer (i.e., server) shares files with other peers.

Cloud-based data stores share more overlaps with Data Grids in the properties listed in Table I. They deliver more abstract storage resources (due to more reliance on virtualization) for different types of workloads in an accurate economic model. They also provide more elastic and scalable resources for different demands in size. These opportunities result in the migration of data-intensive applications to the clouds and cause two categories of issues. One category is more specific to cloud-based data stores and consists of issues such as vendor-lock in, multi-tenancy, network congestion, monetary cost optimization, etc. Another category is common between cloud-based data stores and Data Grids and includes issues like data consistency and latency management. Some of these issues require totally new solutions, and some of them have mature solutions and may be applicable either wholly or with some modifications based on the different properties in cloud-based data stores.

### 2.2. Terms and Definitions

A *data-intensive application* system consists of applications that generate, manipulate, and analyze large amount of data. With the emergence of cloud-based storage services, the data generated from these applications are typically stored in *single data store* or *Geo-replicated data stores* (several data stores in different worldwide locations). The data is organized as *dataset* which is created and accessed across data centers (DCs) by *users/application providers*. *Metadata* describe the dataset with respect to the several attributes such as name, creation time, owner, replicas location, etc.

A dataset consists of a set of *objects* or *records* that are modeled in relational databases (RDBs) or NoSQL databases. The data store that manages RDBs and NoSQL databases are respectively called as *relational data store* and *NoSQL data store*. As shown in Fig. 2, RDBs have fixed and predefined fields for each object whereas NoSQL databases do not. NoSQL data stores use *key-value* data model (or its variations such as graph and document) in which each object is associated with a pair of *key* and *value*. Key is unique and is used to store and retrieve the associated value of an object. NewSQL data stores follow a *hierarchical data model*, which introduces a *child-parent* relation between each pair of table where the child table (Campaign in Fig. 3) borrows the primary key of its parent (e.g., Customer) as a

Table I: Comparison between data-intensive networks in characteristics and objectives.

Property	Cloud-based Data Stores	Data Grids	Content Delivery Network (CDN)	Peer to Peer (P2P)
Purpose	Pay-as-you-go model, on-demand provisioning and elasticity	Analysis, generating, and collaboration over data	File sharing and content distribution	Improving user-perceived latency
Management Entity	Vendor	Virtual Organization	Single organization	Individual
Organization	Tree-based [Wang et al. 2014] <sup>†</sup> Fully optical Hybrid	Hierarchy Federation	Hierarchy	Unstructured Structured Hybrid [Venugopal et al. 2006]
Service Delivery	IaaS, PaaS, and SaaS	IaaS	IaaS	IaaS
Access Type	Read-intensive Write-intensive Equally of both	Read-intensive with rare writes	Read-only	Read-intensive with frequent writes
Data Type	Key-value, Document-based, Extensible record, and Relational	Object-based (Often big chunks)	Object-based (e.g., media, software, script, text)	Object/file-based
Replica Discovery	HTTP requests Replica Catalog	Replica Catalog	HTTP requests	Distributed Hash Table <sup>††</sup> Flooded requests
Replica Placement	See section 4.1.6	Popularity Primary replicas	A primary copy Caching	Popularity without primary replica
Consistency	Weak and Strong	Weak	Strong	Weak
Transaction Support	Only in relational data stores (e.g., Amazon RDS)	None	None	None
Latency Management	Replication, caching, streaming	Replication, caching, streaming	Replication, caching, streaming	Caching, streaming
Cost Optimization	Pay-as-you-go model (in granularity of byte per day for storage and byte for bandwidth)	Generally available for not-for-profit work or project-oriented	Content owners pay CDN operators which, in turn, pays ISPs to host contents	Users Pay P2P to receive sharing files.

<sup>†</sup> Tree-based organization has a flexible topology, while fully optical (consisting of a “pure” optical switching network) and hybrid (including switching network of electrical packet and optical circuit) organizations have a fixed topology. Google and Facebook deploy fat tree topology, a variant of tree topology, in their datacenter architecture. <sup>††</sup> Distributed hash table is used for structured organization and flooded requests for unstructured organization.

prefix for its primary key. In fact, this data model has a *directory table* (e.g., *Customer* in Fig. 3) in the top of hierarchical structure and each row of directory table together with all rows in the descendant tables (e.g., *Campaign* and *AdGroup*) constructs a *directory*.

A *complete* or *partial* replica of the dataset is stored in a single data store or different data stores across Geo-distributed DCs based on the required QoS (response time, availability, durability, monetary cost). An operation to update an object can be initially submitted to a predefined replica (called *single master*) or to any replicas (*multi-master*) based on predefined strategy (e.g., the closest). Replicas are *consistent* when they have the same value for an object. Replicas are in *weak/strong* consistency status if replicas return (probably) *different/same* values for a read operation. A *transaction* is a set of reads and writes, and it is *committed* if all reads and writes are conducted (on all replicas of data); otherwise it is *aborted*.

**Service Level Agreement (SLA)** is a formal commitment between the cloud provider (e.g., Amazon and Microsoft Azure) and the application providers/users. For example, current cloud providers compensate users with service credit if the availability of a data store would be below a specific value in percent.

### 2.3. Data-intensive applications

In respect to the cloud characteristics, two types of data-intensive applications can be nominated for the cloud deployment [Abadi 2009].

*Online transaction processing (OLTP)* applications must guarantee ACID properties and provide an “all-or-nothing” proposition that implies each set of operations in a transaction must complete or no operation should be completed. Deploying OLTP applications across data stores is not straightforward because achieving the ACID properties requires acquiring distributed locks, executing complex commit protocols and transferring data over network, which in turn causes network congestion across data stores and introduces network latency. Thus, OLTP should be adapted to ACID properties at the expense of high latency to serve reads and writes across data stores. *Online Analytical Processing (OLAP)* applications usually use read-only databases and often need to handle complex queries to retrieve the desired data for data

warehouse. The updates in OLAP are conducted on regular basis (e.g., per day or per week), and their rate is lower than that of OLTP. Hence, the OLAP applications do not need to acquire distributed locks and can avoid complex commit protocols.

Both OLTP and OLAP should handle a tremendous volume of data at incredible rates of growth. This volume of data is referred to as *big data* which has challenges in five aspects: *volume* refers to the amount of data; *variety* refers to the different types of generated data; *velocity* refers to the rate of data generation and the requirement of accelerating analysis; *veracity* refers to the certainty of data; and *value* refers to the determination of hidden values from datasets. Schema-less NoSQL databases easily cope with two of these aspects via providing an infinite pool of storage (i.e., volume) in different types of data (i.e., variety). For other aspects, NoSQL data stores are controversial for OLTP and OLAP. From the velocity aspect, NoSQL data stores like BigTable [Chang et al. 2008], PNUTS [Cooper et al. 2008], Dynamo [DeCandia et al. 2007], and Cassandra [Lakshman and Malik 2010] facilitate OLTP with reads and writes in low latency and high availability at the expense of weak consistency. This is a part of the current paper to be discussed. From velocity, veracity and value aspects (the last two aspects are more relevant to OLAP), OLAP requires frameworks like Hadoop,<sup>1</sup> Hive [Thusoo et al. 2010] and Pig<sup>2</sup> as well as algorithms in big data mining to analysis and optimize the complex queries in NoSQL data stores. It is worth to mention these frameworks lack rich query processing on the cloud and change into the data model is a feasible solution which is out of scope of this paper. **In the rest of this section, we focus on goals and challenges of data stores which suit OLTP applications.**

## 2.4. Architecture, Goals, and Challenges of Intra-cloud Storage

This section first describes a layered architecture of data store and then discusses the key goals and challenges of deploying a single data store to manage data-intensive applications.

**2.4.1. Architecture of Intra-Cloud Storage.** Cloud storage architecture shares a storage pool through either a dedicated Storage Area Network (SAN) or Network Attached Storage (NAS).<sup>3</sup> It is composed of a distributed file system, Service Level Agreement (SLA), and interface services. The architecture divides the components by physical and logical functions boundaries and relationships to provide more capabilities [Zeng et al. 2009]. In the layered architecture, each layer is constructed based on the services offered by its underneath layer.

(1) *Hardware layer* consists of distributed storage servers in a cluster that consists of several racks, and each of which has disk-heavy storage nodes. (2) *Storage management* provides services for managing data in the storage media. It consists of the fine-grained services like data replication and erasure coding management, replica recovery, load balancing, consistency, and transaction management. (3) *Metadata management* classifies the metadata of stored data in a global domain (e.g., storage cluster) and collaborates with different domains to locate data when it is stored or retrieved. For example, *Object Table* in Window Azure Storage [Calder et al. 2011] has a primary key containing three properties: *AccountName*, *PartitionName*, and *ObjectName* that determine the owner, location, and name of the table respectively. (4) *Storage overlay* is responsible for storage virtualization that provides data accessibility and is independent of physical address of data. It converts a logical disk address to the physical address by using metadata. (5) *User interface* provides users with primitive operations and allows cloud providers to publish their capabilities, constraints, and service prices in order to help subscribers to discover the appropriate services based on their requirements.

**2.4.2. Goals of Intra-Cloud Storage.** Table II introduces the five main goals of data-intensive applications deployment in a single data store. These goals are as follows:

- **Performance and cost saving.** Efficient utilization of resources has direct influence on cost saving since different data stores offer the pay-as-you-go model. To achieve both these goals, a combination of storage services varying in price and performance yields the desired performance (i.e., response time) with lower cost as compared to relying on one type of storage service. To make an effective combination, the decision on when and which type of storage services to use should be made based on *hot*- and *cold*-spot statuses of data, respectively receiving many and a few read/write requests, and the required QoS. **Current data stores do not guarantee SLA in terms of performance and only allow application providers to measure latency within the storage service (i.e., the time between when a request was served at one of the storage nodes and when the response left the storage node) and the latency over the network (i.e., for the operation to travel from the VM to the storage node and for the response to travel back).**

<sup>1</sup><http://wiki.apache.org/hadoop>

<sup>2</sup><http://hadoop.apache.org/pig/>

<sup>3</sup>A NAS is a single storage device that operates on file system and makes TCP/IP and Ethernet connections. In contrast, a SAN is a local network of multiple devices that operates on disk blocks and uses fiber channel interconnections.

Table II: Intra-Cloud storage goals.

Goals	Techniques (Examples)	Section(s)
Performance and cost saving	Combination of different storage services	§ 7
Fault tolerance and availability	Replication (Random replication, Copyset, MRR)	§ 4.1
	Erasure coding	§ 4.2
Multi-tenancy	Shared table (salesforce.com [Lomet 1996])	—
	Shared process (Pisces [Shue et al. 2012], ElasTras [Das et al. 2013a])	—
	Shared machine (Commercial DaaS, like Amazon RDS)	—
Elasticity and load balancing	Storage-load-aware (Dynamo, BigTable)	—
	Access-aware ([Chen et al. 2013])	—

- *Fault tolerance and availability.* Failures arise from *faulty hardware* and *software* in two models. *Byzantine* model presents arbitrary behaviour and can be survived via  $2f + 1$  replicas ( $f$  is number of failures) along with non-cryptographic hashes or cryptographic primitives. *Crash-stop* model silently stops the execution of nodes and can be tolerated via *data replication* and *erasure coding* in multiple servers located in different racks, which in turn, are in different domains. This failure model has two types: *correlate* and *independent*. The random placement of replicas used in current data stores only tackles independent failures. Copyset [Cidon et al. 2013] replicates an object on a set of storage nodes to tolerate correlated failures, and multi-failure resilient replication (MRR) [Liu and Shen 2016] places replicas of a single chunk in a group of nodes (partitioning into different sets across DCs) to cope with both types of failure.
- *Multi-tenancy.* Multi-tenancy allows users to run applications on shared hardware and software infrastructure so that their isolated performance is guaranteed. It brings benefits for cloud providers to improve infrastructure utilization by eliminating to allocate the infrastructure to the maximum load; consequently cloud providers run applications on and store data in less infrastructure. This, in turn, reduces the monetary cost for users. However, multi-tenancy causes *variable and unpredictable performance*, and *multi-tenant interference and unfairness* [Shue et al. 2012] which happen to different multi-tenancy models (from the weakest to the strongest): *shared table* model in which applications/tenants are stored in the same tables with an identification field for each tenant, *shared process* model in which applications share the database process with an individual table for each tenant, and *shared machine* model in which applications only share the physical hardware, but they have an independent database process and table. Among these models, shared process has the best performance as a database tenant is migrated, while shared machine brings inefficient resource sharing among tenants and shared table is not suitable for applications with different data modes [Das et al. 2013a]. Examples of each model are shown in Table II.
- *Elasticity and load balancing.* Elasticity refers to expansion (scaling up) and consolidation (scaling down) of servers during load changes in a dynamic system. It is orthogonal with load balancing in which workloads are dynamically moved from one server to another under skewed query distribution so that all servers handle workloads almost equally. This improves the infrastructure utilization and even reduces monetary cost. There are two approaches to achieve load balancing. *Storage-load-aware* approach uses *key range* and *consistent hash-algorithm* techniques to distribute data and balance load across storage nodes. This approach, used by almost all existing data stores, is not effective when data-intensive applications experience hot- and cold-spot statuses. These applications thus deploy the *access-load-aware* approach for load balancing [Chen et al. 2013]. Both approaches use either *stop and copy* migration which is economic in data transferring or *live* migration which is network-intensive task but incurs less service interruptions when compared to the former migration technique. Nevertheless, current cloud providers support auto-scaling mechanisms in which the type, maximum and minimum number of database instances should be defined by users. Thus, they cannot provide a fine-grained SLA in this respect.

2.4.3. *Challenges of Intra-Cloud Storage.* Table III introduces what challenges application providers confront with the deployment of their applications within a data store. These challenges are as below:

- *Unavailability of services and data lock-in.* *Data replication* across storage nodes in an efficient way (instead of random replica placement widely used by current data stores) and *erasure coding* are used for high availability of data. Using these schemes across data stores also mitigates data lock-in in the face of appearance of new data store with lower price, mobility of users, and change in workload that demands data migration.
- *Data transfer bottleneck.* This challenge arises when data-intensive applications are deployed within a single data store. It reflects the optimality of data placement that should be conducted based on the characteristics of the application as listed for OLTP and social networks in Table III. To reduce more network congestion, data flows should be monitored in the switches, and they should be then scheduled (i) based on the data flow prediction [Das et al. 2013b], (ii) when data congestion occurs [Al-Fares et al.



Table III: Intra-Cloud storage challenges

Challenges	Solutions	Section(s)/References
Unavailability of services and data lock-in	Data replication Erasure coding	§4.1 §4.2
Data transfer bottleneck	(i) Workload-aware partitioning (for OLTP) (ii) Partitioning of social graph, co-locating the data of a user and his friends along with the topology of DC (for social networks) (iii) Scheduling and monitoring of data flows	[Kumar et al. 2014][Kamal et al. 2016] [Tang et al. 2015][Chen et al. 2016][Zhou et al. 2016] [Al-Fares et al. 2010] [Das et al. 2013b][Shieh et al. 2011]
Performance unpredictability	(i) Data replication and redundant requests (ii) Static and dynamic reservation of bandwidth (iii) Centralized and distributed bandwidth guarantee (iv) Bandwidth guarantee based on network topology and application communications	[Stewart et al. 2013] [Shue et al. 2012] [Ballani et al. 2011][Xie et al. 2012] [Jeyakumar et al. 2013][Ballani et al. 2011][Papa et al. 2013][Guo et al. 2014] [Lee et al. 2014]
Data security	(i) Technical solutions (encryption algorithms, audit third party (ATP), digital signature) (ii) Managerial solutions (iii) A combination of solutions (i) and (ii)	[Shin et al. 2017]

2010], and (iii) for integrated data flows [Shieh et al. 2011] rather than individual data flow. In addition, data aggregation [Costa et al. 2012], novel network typologies [Wang et al. 2014], and optical circuit switching deployment [Chen et al. 2014] are other approaches to drop network congestion.

- *Performance unpredictability.* Shared storage services face the challenges due to multi tenancy, like *unpredictable performance* and *multi-tenant unfairness*, which results in degrading the response time of requests. In Table III, the first two solutions solve challenge with respect to storage, where replicas placement and selection should be considered. The last two solutions cope the challenges related to network aspect, where *fairness* in allocated network bandwidth to cloud applications and *maximizing of network bandwidth utilization* should be taken into consideration. These solutions are: (i) static and dynamic reservation of bandwidth where the static approach cannot efficiently utilize bandwidth, (ii) centralized and distributed bandwidth guarantee where the distributed approach is more scalable, and (iii) bandwidth guarantee based on the network topology, and application communications which is more efficient [Lee et al. 2014]. Such solutions suffer from data delivery within deadline, and they thus are suitable for batch applications, but not for OLTP applications which require the completion of data flows within deadline [Vamanan et al. 2012].
- *Data security.* This is one of the strongest barriers in the adoption of public clouds to store users' data. It is a combination of (i) data integrity, protecting data from any unauthorized operations; (ii) data confidentiality, keeping data secret in the storage (iii) data availability, using data at any time and place; (iv) data privacy, allowing data owner to selectively reveal their data; (v) data transition, detecting data leakage and lost during data transfer into the storage; and (vi) data location, specifying who has jurisdiction and legislation over data in a transparent way. Many solutions were proposed in recent decades to deal with concerns over different security aspects except *data location* where data is probably stored out of users' control. These solutions are not generally applicable since unlike the traditional systems with two parties, the cloud environment includes 3 parties: users, storage services, and vendors. The last party is a potential threat to the security of data since it can provide a secondary usage for itself (e.g., advertisement purposes) or for governments.

Solutions to relieve concerns over data security in the cloud context can be *technical*, *managerial*, or *both*. In addition to technical solutions as listed in Table III, managerial solutions should be considered to relieve security concerns relating to the geographical location of data. A combination of both type of solutions can be: (i) designing location-aware algorithms for data placement as data are replicated to reduce latency and monetary cost, (ii) providing location-proof mechanisms for user to know the precise location of data (e.g., measuring communication latency and determining distance), and (iii) specifying privacy acts and legislation over data in a transparent way for users. Since these acts, legislation, and security and privacy requirements of users are a fuzzy concept, it would be relevant to design a fuzzy framework in conjunction of location-aware algorithms. This helps users to find storage services which are more secure in respect to rules applied by the location of data.

A better approach to achieve the discussed goals and avoid the above challenges is to store data across data stores. In the rest of the section, we discuss the benefits of this solution as well as the challenges that arise from it.

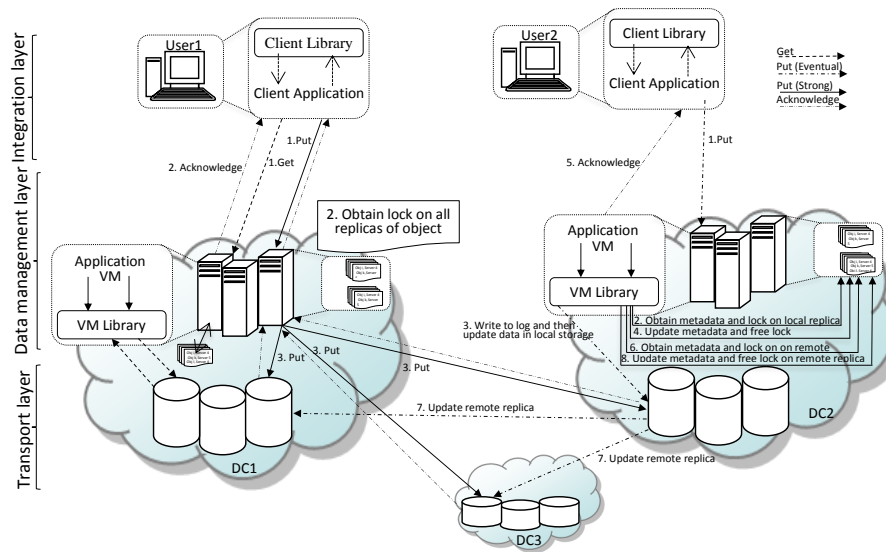




Table IV: Inter-Cloud storage goals.

Goals	Techniques/Schemes	Section(s)
High availability, durability, and data lock-in avoidance	Data replication	§4.1
	Erasure coding	§4.2
Cost benefit	Exploitation of pricing differences across data stores and time-varying workloads	§7
Low user-perceived latency	Placing replicas close to users	—
	Co-locating the data accessed by the same transactions	
	Determining the location and roles of replicas (master and slave) in a quorum-based configuration [Sharov et al. 2015]	
	Cryptographic protocols with erasure coding and RAID techniques	
High data confidentiality, integrity, and auditability	Cryptographic protocols with erasure coding and RAID techniques	—

2.5.2. *Goals of Inter-Cloud Storage.* Table IV lists the key goals of leveraging Geo-replicated data stores. These are as follows:

- *High availability, durability, and data lock-in avoidance.* These are achievable via data replication scheme across data stores owned by different cloud providers. A 3-way data replication provides availability of 7 nines [Mansouri et al. 2013] which is adequate for most applications. This is inadequate for protecting data against *correlated* failures, while two replicas are sufficient for guarding data against *independent* failures [Cidon et al. 2015]. *Erasure coding* is another way to attain higher durability even though it degrades data availability in comparison to data replication.
- *Cost benefit.* Due to *pricing differences* across data stores and *time-varying workloads*, application providers can diversify their infrastructure in terms of vendors and locations to optimize cost. The cost optimization should be integrated with the demanding QoS level like availability, response time, consistency level, etc. This can be led to a trade-off between the cost of two resources (e.g., storage vs. computing) or the total cost optimization based on a single-/multi-QoS metrics.
- *Low user-perceived latency.* Application providers achieve lower latency through deploying applications across multiple cloud services rather than within a single cloud services. Applications can further decrease latency via techniques listed in Table IV. In spite of these efforts, users may still observe the latency variation which can be improved through caching data in memory [Nishtala et al. 2013]), issuing redundant reads/writes to replicas [Wu et al. 2015a]), and using feedback from servers and users to prevent requests redirection to saturated servers [Suresh et al. 2015]).
- *High data confidentiality, integrity, and auditability.* Using *Cryptographic protocols with erasure coding and RAID techniques* on top of multiple data stores improves security in some aspects as deployed in HAIL [Bowers et al. 2009] and DepSky [Bessani et al. 2011]. In such techniques, several concerns are important: scalability, cost of computation and storage for encoding data, and the decision on where the data is encoded and the keys used for data encryption are maintained. A combination of private and public data stores and applying these techniques across public data stores offered by different vendors improve data protection against both insider and outsider attackers, especially for insider ones who require access to data in different data stores. Data placement in multiple data stores, on the other hand, brings a side effect since different replicas are probably stored under different privacy rules. Selection of data stores with similar privacy acts and legislation rules would be relevant to alleviate this side effect.

2.5.3. *Challenges of Inter-Cloud Storage.* The deployment of data-intensive applications across data stores is faced with the key challenges as listed in Table V. These are discussed as below:

- *Cloud interoperability and portability.* Cloud interoperability refers to the ability of different cloud providers to communicate with each other and agree on the data types, SLAs, etc. Cloud portability means the ability to migrate application components and data across cloud providers regardless of APIs, data types, and data models. Table V lists solutions for this challenge.
- *Network congestion.* Operating across Geo-DCs causes *network congestion* which can be *time-sensitive* or *non time-sensitive*. The former, like interactive traffic, is sensitive to delay, while the latter, like transferring big data and backing up data, is not so strict to delay and can be handled within deadline [Zhang et al. 2015] or without deadline [Jain et al. 2013]. The first solution for this challenge, listed in Table V, is expensive, while the second solution increases the utilization of network. The last solution is Software-Defined Networking (SDN) [Kreutz et al. 2015]. SDN separates *control plane* that decides how to handle network traffic, and *data paths* that forwards traffic based on the decision made from control plane. All these solutions answer a part of this fundamental question: *how to schedule the data transfer so that it is completed within a deadline and budget subject to the guaranteed network fairness and throughput for jobs that processes the data.*
- *Strong consistency and transaction guarantee.* Due to high communication latency between DCs, coordination across replicas to guarantee strong consistency can drive users away. To avoid high communication latency, some data stores compromise strong consistency at the expense of application semantics

Table V: Inter-Cloud storage challenges

Challenges	Solutions	Example	Section(s)
Cloud interoperability and portability	Standard protocols for IaaS Abstraction storage layer Open API	n/a† CASL[Hill and Humphrey 2010] JCloud††	—
Network congestion	Dedicating redundant links across DCs Store and forward approach Software-Defined Networking (SDN)	n/a Postcard [Feng et al. 2012], [Wu et al. 2015b] B4 [Jain et al. 2013], [Wu et al. 2015b]	—
Strong consistency and transaction guarantee	Heavy-weight coordination protocols Contemporary techniques	See Appendixes C, D, and E	§5, §6

†n/a: not applicable, ††JCloud:<https://jclouds.apache.org/>

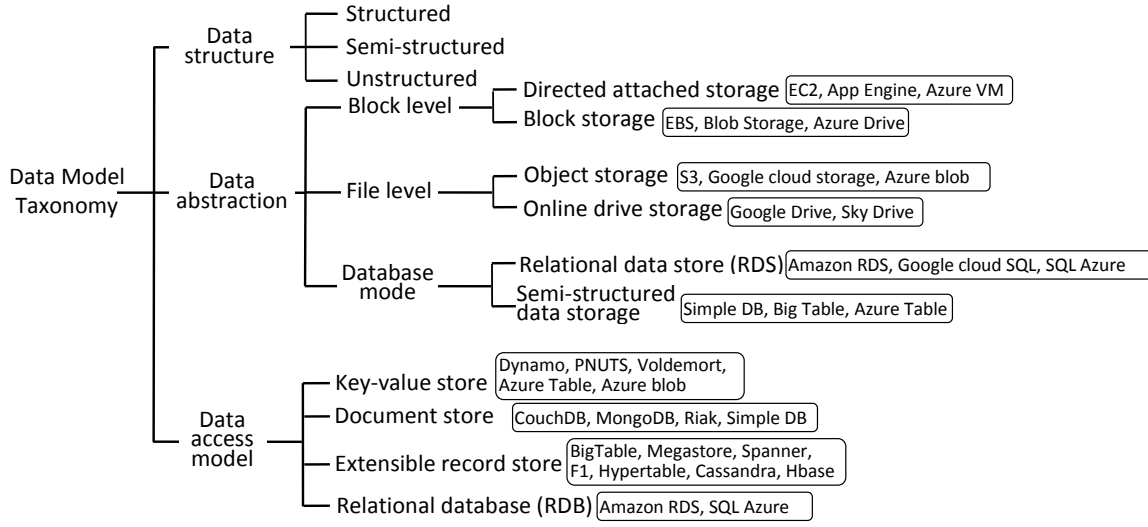


Fig. 5: Data model taxonomy

violations and stale data observations. Others, on the other hand, provide strong consistency in the cost of low availability. To achieve strong consistency without compromising with availability and scalability, coordination across replicas should be reduced or even eliminated.

Our survey paper focuses on some of the goals and challenges discussed above. In particular, the rest of the paper discusses the five elements of data storage management specified in Fig. 1.

### 3. DATA MODEL

Data model reflects that how data is logically organized, stored, retrieved, and updated in data stores. We thus study it from different aspects and map data stores to the provided data model taxonomy in Fig. 5.

#### 3.1. Data structure

Data structure affects the speed of assimilation and information retrieving. It has three categories. (i) *Structured* data refers to data that defines the relation between the fields specified with a name and value, e.g., RDBs. It supports a comprehensive query and transaction processing facilities. (ii) *Semi-structured* is associated with the special form of structured data with a specific schema known for its application and database deployments (e.g., document and extensible DBs). It supports simple query (e.g., primitive operations) and transaction facilities as compared to structured data. (iii) *Unstructured* data refers to data that have neither pre-defined data model nor organized in a pre-defined way (e.g., video, audio, and heavy-text files). It takes the simplest data access model, i.e., key-value, that delivers high scalability and performance at the cost of sacrificing data consistency semantic and transaction support.

In these logical data structures, data is internally organized row-by-row, column-by-column closely related to database normalization, or combination of both schemes- called hybrid – within a block. *Structured* data can be organized in all schemes, *semi-structured* in row-by-row and hybrid schemes, and *Unstructured* data in a row-by-row scheme.

#### 3.2. Data abstraction

This refers to different levels of storage abstraction in data stores. These levels are as below:

Table VI: Comparison between different storage abstractions.

Ease of use†	Scalability	Performance††	Cost [Applicability]
File-level	File-level	Block-level	Block-level [OLAP applications]
Database mode	Block-level	Database mode	Database mode [OLTP applications with low latency queries]
Block-level	Database-mode	File-level	File-level [Backup data and web content static]

† The levels of storage abstract are listed from high to low for each aspect listed in each column. †† Performance is defined in terms of accessibility.

- (1) *Block-level* provides the fastest access to data for virtual machine (VM). It is classified into (i) *directed-attached storage* coming with a VM that provides highly sequential I/O performance with a low cost, and (ii) *block storage* that pairs with a VM as a local disk.
- (2) *File-level* is the most common form of data abstraction due to its ease of management via simple API. It is provided in the forms of (i) *object storage* that enables users to store large binary objects anywhere and anytime, and (ii) *online drive storage* that provides users with folders on the Internet as storage.
- (3) *Database mode* offers storage in the forms of *relational data store* and *semi-structured data storage* which respectively provide users with RDB and NoSQL/NewSQL databases. RDB exploiting the SQL standard does not scale easily to serve large web applications but guarantees strong consistency. In contrast, NoSQL provides horizontal scalability by means of shared nothing, replicating, and partitioning data over many servers for simple operations. In fact, it preserves BASE (Basically Available, Soft state, Eventually consistent) properties instead of ACID ones in order to achieve higher performance and scalability. NewSQL— as a combination of *RDB* and *NoSQL*— targets delivering the scalability similar to NoSQL, meanwhile maintaining ACID properties.

Table VI compares different levels of storage abstractions in several aspects as well as their applicability. This comparison indicates that as the storage abstraction (ease of use) level increases, the cost and performance of storage reduce.

### 3.3. Data access model

This reflects storing and accessing model of data that affect on consistency and transaction management. It has four categories as below:

- (1) *Key-value database* stores keys and values which are indexed by keys. It supports primitive operations and high scalability via keys distribution over servers. Data can be retrieved from the data store based on more than one attribute if additional key-value indexes are maintained.
- (2) *Document database* stores all kinds of documents indexed in the traditional sense and provides primitive operations without ACID guarantee. It thus supports *eventual consistency* and achieves scalability via *asynchronous* replication, *shard* (i.e., horizontal partition of data in the database), or both.
- (3) *Extensible record database* is analogous to table in which columns are grouped, and rows are split and distributed on storage nodes [Cattell 2011] based on the primary key range as a *tablet* representing the unit of distribution and load balancing. Each cell of the table contains multiple versions of the same data that are indexed in decreasing timestamps order, thereby the most recent version can always read first [Sakr et al. 2011]. This scheme is called NewSQL which is equivalent with *entity group* in Megastore, *shard* in Spanner, and *directory* in F1 [Shute et al. 2013].
- (4) *Relational database* (RDB) has a comprehensive pre-defined scheme and provides manipulation of data through SQL interface that supports ACID properties. Except for *small-scope* transactions, RDB cannot scale the same as NoSQL.

Table VII: Comparison between NoSQL (the first three databases) and Relational databases.

Database	Simplicity	Flexibility	Scalability	Properties	Data Structure	SLA	Application [Query] type
key-value	High	High	High	—	(Un-/)structured	Weak	OLAP [Simple]
Document	High	Moderate	Moderate	BASE	Semi-structured	Weak	OLAP [Moderate]
Extensible record	High	High	High	ACID	(Semi-/)structured	Weak	OLTP [Moderate, repetitive]
Relational	Low	Low	Low	ACID	Structured	Weak	OLTP [Complex]

Table VII compares different databases in several aspects and indicates which type of application and query can deploy them. NoSQL databases offer horizontal scalability and high availability by sacrificing *consistency semantic* which makes them unsuitable for the deployment of OLTP applications. To provide stronger consistency semantic for OLTP applications, it is vital to carefully partition data within and especially across data stores and to use the mechanisms that exempt or minimize the coordination across transactions. Moreover, these databases have several limitations relating to big data for OLAP applications. All these disk-based data stores cannot suitably facilitate OLAP applications in the concept of *velocity* and thus most commercial vendors combine them with in-memory NoSQL/relational data stores

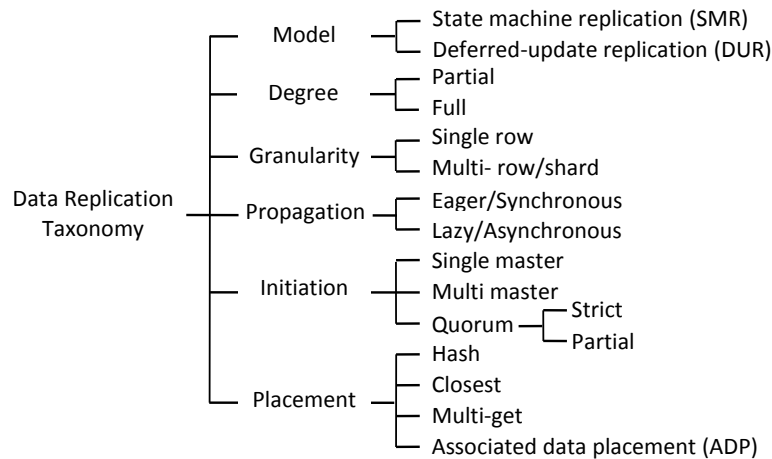


Fig. 6: Data replication taxonomy

(e.g., Memcached,<sup>5</sup> Redis,<sup>6</sup> and RAMCloud [Rumble et al. 2014]) to further improve performance. These databases also cannot be completely fitted with the concept of *velocity* relating to big data for OLAP applications which receive data with different formats. It is thus required a platform to translate these formats into a canonical format. OLAP applications can combat the remaining limitations (i.e., veracity and value) arising from using these databases in the face of large data volumes via designing indexes to retrieve data. Moreover, NoSQL data stores guarantee weak SLA in terms of availability [Sakr 2014] and auto-scaling, and without any SLA in the case of response time. In fact, they measure response time as data is stored and retrieved, and they also replicate data across geographical locations on users' request.

Other class of database is relational. This class of database compromise high availability and scalability for stronger consistency semantics and for providing higher query processing facilities. Current cloud providers offer this class of database in the form of Database-as-a-Service (DaaS) such as Amazon RDS, SQL Azure, and Google cloud SQL. DaaS allow application providers to use the full capabilities of MySQL, MariaDB, PostgreSQL, Oracle, and Microsoft SQL Server DB engines. In fact, DaaS support complex queries, distributed transactions, store procedures, views, etc., while they confront with several limitations mainly elasticity, strong consistency, transaction management across geographical locations, and live migration for purposes like multi-tenancy and elasticity. Moreover, similar to NoSQL databases, DaaS support almost the same SLA and only can scale out/down based on manual rules defined by users according to the workload. DaaS require more complex mechanisms to guarantee SLA because they should manage resources like CPU, RAM, and disk. For more details on challenges and opportunities, the readers are referred to the survey on cloud-hosted databases proposed by Sakr [Sakr 2014].

Fig. 5 provides examples of discussed classes of databases, and Table I in the Appendix summarizes them in several main aspects. For more details on these databases, readers are referred to [Sakr et al. 2011].

## 4. DATA DISPERSION

This section discusses data dispersion schemes as shown in Fig. 1.

### 4.1. Data replication

Data replication improves availability, performance (via serving requests by different replicas), and user-perceived latency (by assigning requests to the closest replica) at the cost of replicas coordination and storage overheads. This is affected by facets of data replication based on the taxonomy in Fig. 6.

**4.1.1. Data replication model.** There are two replication models for fault tolerant data stores [Pedone et al. 2000]: The model is *state machine replication* (SMR) in which all replicas receive and process all operations in a *deterministic way* (in the same order) using *atomic broadcast* (see §6.1.3). This implies SMR is abort-free and failure transparency which means if a replica fails to process some operations those are still processed in other replicas. However, SMR has low throughput and scalability for read and write (RW) transactions since all servers process each transaction. Thus, the scalability and throughput are confined by the processing power of a server. Scalable-SMR (S-SMR) [Bezerra et al. 2014] solves this issue across data stores via (i) partitioning database and replicating each partition, and (ii) using cache techniques

<sup>5</sup><https://memcached.org/>

<sup>6</sup><https://redis.io/>

Table VIII: Comparison between replication models.

Replication Models	Advantages	Disadvantages	Consistency semantic
State machine replication <sup>†</sup>	Failure transparency Abort-free	Low throughput and scalability for RW transactions	Linearizability
Deferred-update replication	High throughput and scalability for RO transactions	Stale data Replica divergence	non-Linearizability

<sup>†</sup> It is also called *active replication*.

Table IX: Comparison between full and partial replications.

Properties	Comparison	Reason/Description
Scalability	Partial > Full	Due to access to a subset of DCs not all DCs
Complexity	Partial > Full	Due to the requirement for the exact knowledge where data reside
Storage Cost	Partial < Full	Due to data replication in a subset of data stores
Applicability	Partial > Full	If read requests come a specific DCs, or when objects are write-intensive
	Partial < Full	If read requests come from all DCs, or when transactions are global

to reduce communication between partitions without compromising consistency. SMR and S-SMR are suitable for contention-intensive and *irrevocable transactions* that require abort-free execution.

The second model is Deferred-update replication (DUR). It resembles single-/multi master replication and scales better than SMR due to locally executing RW transactions on a server and then propagating updates to other servers. Thus, in DUR, the RW transactions do not scale with the number of replicas in comparison to the read-only (RO) transactions executed only on a server without communication across servers by using a *multiversion* technique. Scalable-DUR (S-DUR) [Sciascia et al. 2012] and Parallel-DUR (P-DUR) [Pacheco et al. 2014] allow update transactions to scale with the number of servers and the number of cores available for a replica respectively. In respect to pros and cons summarized in Table VIII, the scalability and throughput of transactions can be improved through borrowing the parallelism in DUR and abort-free feature in SMR [Kobus et al. 2013].

**4.1.2. Data replication degree.** Data replication can be either *partial* or *full*. In partial (resp. full) replication each node hosts a *portion* (resp. *all*) of data items. For example, in the context of Geo-DCs, *partial* (resp. *full*) replication means that certain (resp. all) DCs contain a replica of certain (resp. all) data items.

Table IX shows that partial replication outperforms full replication in storage services cost due to access to a subset of data stores deployed across DCs. It is also better than full replication in scalability because full replication is restricted by the capacity of each replica that certifies and processes transactions. These advantages demand more complex mechanisms for consistency and transaction support in partial replication, which potentially degrades response time. Many partial replication protocols provide such mechanisms at the expense of communication cost same as full replication [Armendáriz-Inigo et al. 2008]. This is due to unpredictable *overlapping* transactions in which the start time of transaction  $T_i$  is less than the commit time of transaction  $T_j$  and the intersection of write set  $T_i$  and  $T_j$  is not empty. The deployment of *genuine partial replication* solves such issue and enforces a transaction to involve only the subset of servers/DCs containing the desired replicas for coordination. In terms of applicability, partial and full replication is more suitable for write-intensive objects (due to submitting each request to a subset of DCs [Shen et al. 2015]) and for execution of global (multi-shard) transactions respectively.

Therefore, the characteristics of workload and the number of DCs are main factors in making a decision on what *data replication degree* should be selected. If the number of DCs is small, full replication is preferable; otherwise, if global transactions access few DCs, partial replication is a better choice.

**4.1.3. Data replication granularity.** This defines the level of data unit that is replicated, manipulated, and synchronized in data stores. Replication granularity has two types: *single row* and *multi-row/shard*. The former naturally provides horizontal data partitioning, thereby allowing high availability and scalability in data store like Bigtable, PNUTS, and Dynamo. The latter is the first step beyond single row granularity for new generation web-applications that require to attain both high scalability of NoSQL and ACID properties of RDBs (i.e., NewSQL features). This type of granularity has an essential effect on the scalability of transactions, and according to it, we classify transactions from granularity aspect in §6.1.5.

**4.1.4. Update propagation.** This reflects *when* updates take place and is either *eager/synchronous* or *lazy/asynchronous*. In eager propagation, the committed data is simultaneously conducted on all replicas, while in lazy propagation the changes are first applied on master replica and then on slave replicas. Eager propagation is applicable on a single data store like SQL Azure [Campbell et al. 2010] and Amazon RDS, but it is hardly feasible across data stores due to response time degradation and network bottleneck. In contrast, lazy propagation is widely used across data stores (e.g., Cassandra) to improve response time.



Table X: Comparison between update initiation.

Update initiation	Advantages	Disadvantages	Data store(s)
Single Master†	Strong consistency†† Conflict avoidance	Low throughput of RW transactions Single point of failure and bottleneck	Amazon RDS◊ PNUT
Multi-master	High throughput of RO and RW transactions, not performance bottleneck	Stale data Replica divergence	Window Azure Storage (WAS)
Quorum	Adaptable protocol for availability vs. consistency	Determining the location of coordinator Determining the read and write quorum	Cassandra, Riak, Dynamo, Voldemort

†It is also called *primary backup*. ††Strong consistency is provided with *eager propagation*, not lazy. ◊ Amazon RDS: <https://aws.amazon.com/rds/>

**4.1.5. Update initiation.** This refers to *where* updates are executed in the first place. Three approaches for update initiation are discussed as below:

*Single master* approach deploys a replica at the closest distance to the user or a replica receiving the most updates as master replica. All updates are first submitted to the master replica and then are propagated either eagerly or lazily to other replicas. In single master with lazy propagation, replicas receive the updates in the same order accepted in the master and might miss the latest versions of updates until the next re-propagation by the master. Single master approach has advantages and disadvantages as listed in Table X. These issues can be mitigated somehow by *multi-master* approach in which every replica can accept update operations for processing and in turn propagates the updated data to other replicas either eagerly or lazily. Thus, this approach increases the throughput of read and write transactions at the cost of stale data, while replicas might receive the updates in the different order which results in replicas divergence and thus the need for conflict resolution.

*Quorum* approach provides a protocol for availability vs. consistency in which writes are sent to a write set/quorum of replicas and reads are submitted to a read quorum of replicas. The set of read quorum and write quorum can be different and both sets share a replica as coordinator/leader. The reads and writes are submitted to a *coordinator replica* which is a single master or multi-master. This protocol suffers from the disadvantages in determining the coordinator location and the quorum of write and read replicas as addressed when workload changes [Sharov et al. 2015]. Though this classical approach guarantees strong consistency, many Geo-replicated data stores achieve higher availability at a cost of weaker consistency via its adaptable version in which a read/write is sent to all replicas and is considered successful if the acknowledgements are received from a quorum (i.e., majority) of replicas. This adapted protocol is configured with write ( $W$ ) and read quorum ( $R$ ) in synchronous writes and reads. The configuration is determined in (i) *strict quorum* in which any two quorums have non-empty intersection (i.e.,  $W + R > N$ , where  $N$  is the number of replicas) to provide strong consistency, and (ii) *partial quorum* in which at least two quorums should not overlap (i.e.,  $W + R < N$ ) to support weak consistency. Generally speaking, (i) a raise in  $\frac{W}{R}$  improves the consistency, and (ii) a raise in  $W$  reduces availability and increases durability.

**4.1.6. Replica placement.** This is related to the mechanism of replica placement in data store and is composed of four categories. (1) *Hash mechanism* is intuitively understood as a random placement and determines the placement of objects based on the hashing outcome (e.g., Cassandra). It effectively balances the load in the system, however it is not effective for a transactional data store that requires co-located multiple data items. (2) *Closest mechanism* replicates a data item in the node which receives the most requests for this data item. Although closest mechanism decreases the traffic in the system, it is not efficient for a transactional data store because the related data accessed by a transaction might be placed in different locations. (3) *Multiget* [Nishtala et al. 2013] seeks to place the associated data in the same location without considering the localized data serving. (4) *Associated data placement* (ADP) [Yu and Pan 2015] makes the strike between *closest* and *multiget* mechanisms.

As discussed above, using each strategy differed on various aspects of replication can affect the latency, consistency and transaction complexity, and even monetary cost. Among these, the key challenge is how to make a trade-off between consistency and latency, where update *initiation* and *propagation* are main factors. **Data stores provide lesser latency and better scalability as these factors mandate weaker consistency (Table II -Appendix).**

## 4.2. Erasure coding

Cloud file system uses *erasure coding* to reduce storage cost and to improve availability as compared to data replication. A  $(k, m)$ -erasure coding divides an object into  $k$  equally sized chunks that hold original data along with  $m$  extra chunks that contain data coding (parity). All these chunks are stored into  $n = k + m$  disks which increases the storage overhead by a factor of  $1/r = k/n < 1$  and tolerates  $m$  faults, as opposed to  $m - 1$  faults for  $m$ -way data replication with a factor of  $m - 1$  storage overhead. For example, a  $(3,2)$ -erasure coding tolerates 2 failed replicas with a  $2/3(=66\%)$  storage overhead as compared to 3-

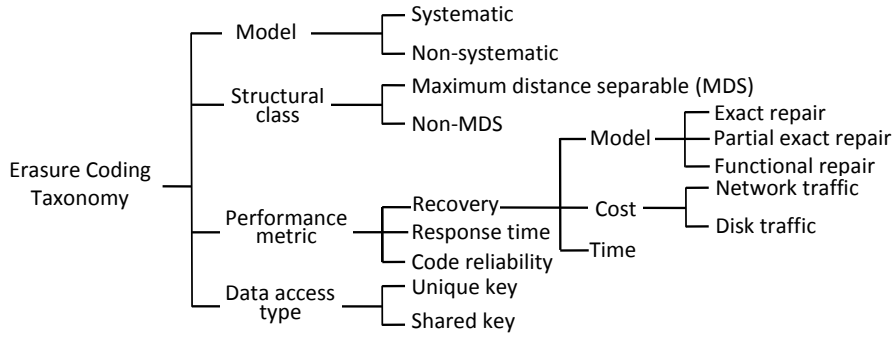


Fig. 7: Erasure coding taxonomy

Table XI: Comparison between Replication, Reed-Solomon Code, and Local Reconstruction Code.

Schemes	Structural Class	Recovery Cost <sup>†</sup>	Storage Overhead	Applicability
(m-way) Replication	n/a	1	(m-1)X	3-way replication in WAS, S3, and Google storage
(k,m)-RS	MDS	K	$(1 + \frac{m}{k})X$	(6,3)-RS in Google and (10,4)-RS in Facebook
(k,l,r)-LRC	Non-MDS	$[r+1, r+1]$	$(1 + \frac{k+r}{k})X$	(12,2,2)-LRC in WAS [Calder et al. 2011]

<sup>†</sup> n/a: not applicable, *Recovery time*: (m-way) Replication < (k,l,r)-LRC < (k,m)-RS. *Durability*: (m-way) Replication < (k,m)-RS < (k,l,r)-LRC.

way replication with the same fault tolerance and a 200% storage overhead. To use erasure coding as an alternative to data replication, we need to investigate it in the following aspects as shown in Fig. 7.

**4.2.1. Model.** Erasure coding has two models. *Systematic* model consists of (*original*) data and *parity* chunks which are separately stored in  $n - m$  and  $m$  storage nodes. *Non-systematic* model includes the coded data (not original data) which are stored in  $n$  nodes. Systematic and non-systematic codes are respectively relevant for archival and data-intensive applications due to respectively low and high rate of reads. Systematic codes seem more suitable for data stores because they decode data when a portion of data is unavailable. Comparatively, non-systematic codes decode data whenever data are retrieved due to storing the coded data not the original data. This may degrade the response time.

**4.2.2. Structural class.** This represents the reconstruct-ability of code that is largely reliant on erasure coding rate (i.e.,  $r$ ), including two classes. The first class is (k,m)-maximum distance separable (MDS) code in which the data is equally divided into  $k$  chunks which are stored in  $n = k + m$  storage nodes. A (k,m)-MDS code tolerates any  $m$  of  $n$  failed nodes and rebuilds the unavailable data from any  $k$  surviving nodes, known as *MDS-property*. A MDS code is optimal in terms of reliability vs. storage overhead while it incurs significant recovery overheads to repair failed nodes. The second class is non-MDS code. Any code that is not MDS code is called non-MDS code. This code tolerates less than  $m$  failed nodes and typically rebuilds the failed nodes from less than  $k$  surviving nodes. Thus, compared to MDS codes, non-MDS codes are (i) more economical in network cost to rebuild the failed nodes, which makes them more suitable for deploying across data stores and (ii) less efficient in the storage cost and fault-tolerance.

**4.2.3. Performance metrics.** Erasure coding is evaluated based on the following performance metrics that have received significant attention in the context of cloud.

(1) *Recovery* refers to the amount of data retrieving from disk to rebuild a failed data chunk. Recovery is important in the below aspects.

- *Recovery model* in which a failed storage node is recovered through survivor nodes has three models [Suh and Ramchandran 2011]. The first is *exact-repair* in which the failed nodes are exactly recovered, thus lost parity with their exact original data are restored. The second is *partial exact-repair* in which the data nodes are fixed exactly and parity nodes are repaired in a functional manner by using random-network-coding framework [Dimakis et al. 2010]. This framework allows to repair a node via retrieving functions of stored data instead of subset of stored data so that the code property (e.g., MDS-property) is maintained. The third is *functional repair* in which the recovered nodes contain different data from that of the failed nodes while the recovered system preserves the MDS-code property. Workload characteristics and the deployed code model determine which recovery model satisfies the application requirements. Exact- and partial exact-repair are appropriate for archival applications while functional repair is suitable for non-secure sensitive applications because it requires the dynamics of repairing and decoding rules that results in information leakage.
- *Recovery cost* is the total number of required chunks to rebuild a failed data chunk. It consists of *disk traffic* and *network traffic* affected by network topology and replication policy [Zhang et al. 2010]. Re-

cently, the trade-off between recovery cost and storage cost takes considerable attention in the context of cloud as discussed below.

A  $(k, m)$ -Reed-Solomon (RS) code contains  $k$  data chunks and  $m$  parity chunks, where each parity chunk is computed from  $k$  chunks. When a data chunk is unavailable, there is always a need of any subset of  $k$  chunks from  $m + k$  chunks, as recovery cost, to rebuild the data chunk. This code is used in Google ColossusFS [Ford et al. 2010] and Facebook HDFS [Muralidhar et al. 2014] within a DC (a XOR-based code—using pure XOR operation during coding computation— across DCs). In spite of the RS code optimality in reliability vs. storage overhead, it is still unprofitable due to high bandwidth requirements within and across data stores. Hitchhiker [Rashmi et al. 2014] mitigates this issue without compromise on storage cost and fault tolerance throughout the adapted RS code in which a single strip RS code is divided into two correlated sub-stripes.

Similar to MDS-code, *regenerating* and non-MDS codes [Wu et al. 2007] alleviate the network and disk traffics. Regenerating codes aim at the optimal trade-off between storage and recovery cost and come with two optimal options [Rashmi et al. 2011]. The first is the *minimum storage regenerating* (MSR) codes which minimize the recovery cost keeping the storage overheads the same as that in MDS codes. NCCloud [Chen et al. 2014] uses *functional* MSR, as a non-systematic code, and maintains the same fault tolerance and storage overhead as in RAID-6. It also lowers recovery cost when data migrations happen across data stores due to either transition or permanent failures. The second is the *minimum bandwidth regenerating* (MBR) codes that further minimize the recovery cost since they allow each node to store more data.

A  $(k, l, r)$ —Local Reconstruction Code (LRC) [Huang et al. 2012] divides  $k$  data blocks into  $l$  local groups and creates a *local parity* for each local group and  $r = \frac{k}{l}$  *global parities*. The number of failure it can tolerate is between  $r + 1$  and  $r + l$ . HDFS-Xorbas [Sathiamoorthy et al. 2013] exploits LRC to make a reduction of 2x in the recovery cost at the expense of 14% more storage cost. HACFS [Xia et al. 2015] code also uses LRC to provide a *fast code* with low recovery cost for *hot* data and exploits Product Code (PC) [Roth 2006] to offer a compact code with low storage cost for *cold* data. Table XII compares RS-code and LRC, used in current data stores, with data replication in the main performance metrics.

- *Recovery time* refers to the amount of time to read data from disk and transfer it within and across data stores during recovery. As recovery time increases, *response time* grows, resulting in notorious effect on the data availability. Erasure codes can improve recovery time through two approaches. First, erasure codes should reduce network and disk traffics to which RS codes are inefficient as they read all data blocks for recovery. However, *rotated* RS [Khan et al. 2012] codes are effective due to reading the requested data only. LRC and XOR-base [Rashmi et al. 2014] codes are also viable solutions to decrease recovery time. Second, erasure codes should avoid retrieving data from hot nodes for recovery by replicating hot nodes' data to cold nodes or caching those data in dynamic RAM or solid-state drive.

(2) *Response time* indicates the delay of reading (resp. writing) data from (resp. in) data store and can be improved through the following methods. (i) *Redundant requests* simultaneously read (resp. write)  $n$  coded chunks to retrieve (resp. store)  $k$  parity chunks [Shah et al. 2014]. (ii) *Adaptive batching* of requests makes a trade-off between delay and throughput, as exploited by S3 (Simple Storage Service) and WAS [McCullough et al. 2010]. (iii) *Deploying erasure codes across multiple data stores* improves availability and reduces latency. Fast Cloud [Liang and Kozat 2014] and TOFEC [Liang and Kozat 2015] use the first two methods to make a trade-off between throughput and delay in key-value data stores as workload changes dynamically. Response time metric is orthogonal to data monetary cost optimization and is dependent on three factors: scheduling read/write requests, the location of chunks, and parameters that determine the number of data chunks. Xiang et al. [2014] considered these factors and investigated a trade-off between latency and storage cost within a data store.

(3) *Reliability* indicates the *mean time to data loss* (MTTDL). It is estimated by standard Markov model and is influenced by the speed of block recovery [Sathiamoorthy et al. 2013] and the number of failed blocks that can be handled before data loss. LRCs are better in reliability than RS codes, which in turn, are more reliable than replication [Huang et al. 2012] with the same failure tolerance. As already discussed, failures in data stores can be *independent* or *correlated*. To relieve the later failure, the parity chunks should be placed in different racks located in different domains.

**4.2.4. Data access type.** There are two approaches to access chunks of the coded data: *unique key* and *shared key* [Liang and Kozat 2015], in which a key is allocated to a chunk of coded data and the whole coded data respectively. These approaches can be compared in three aspects. (1) *Storage cost*: both approaches are almost the same in the storage cost for writing into a file. In contrast, for reading chunks, shared key is more cost-effective than unique key. (2) *Diversity in delay*: with unique key, each chunk, treated as an individual object, can be replicated in different storage units (i.e., server, rack, and DC). With shared key, chunks are combined into an object and very likely stored in the same storage unit. Thus, in unique (resp. shared) key, there is low (resp. high) correlation in the access delay for different

chunks. (3) *Universal support*: unique key is supported by all data stores, while shared key requires advanced APIs with the capability of partial reads and writes (e.g., S3).

#### 4.3. Hybrid scheme

Hybrid scheme is a combination of data replication and erasure coding schemes to retain the advantages of these schemes while avoiding their disadvantages for data redundancy within and across data stores. Table XII compares two common schemes in performance metrics to which three factors contribute into when and which scheme should be deployed: Access rate (AR) to objects, object size (OS), price (Pr) and performance (Pe) of data stores.

These factors have a significant effect on the *storage overhead*, *recovery cost*, and *read/write latency*. As indicated in Table XII, replication incurs storage overhead more than erasure coding especially for large objects, while it requires less recovery cost due to retrieving the replica from a single server/data store instead of fragmented objects from multiple servers/data stores. Erasure coding is more profitable in read/write latency (i) for cold-spot objects since update operations require re-coding the whole object, and (ii) for large objects because the fragmented objects can be accessed in parallel from multiple server/data stores. For the same reasons, replication is more efficient for hot-spot objects with small size like metadata objects. Thus, cold-spot objects with large size should be distributed across cost-effective data stores in the form of erasure coding, and hot-spot objects with small size across performance-efficient data stores in the form of replication.

We classify the hybrid scheme into two categories. (i) *Simple hybrid* stores an object in the form of either replication or erasure coding (ROC) or replication and erasure coding (RAC) during its lifetime. (ii) *Replicated erasure coding* contains replicas of each chunk of coded objects and its common form is *double coding* which stores two replicas of each coded chunk of object. Compared to ROC, double coding and RAC increase storage cost two times, but they are better in availability and bandwidth cost due to retrieving the lost chunk of the object from the server which has a redundant copy. Table XIII summarizes projects using common redundancy or hybrid schemes. Neither workload characteristics nor data stores diversity (in performance and cost) are fully deployed in these projects using hybrid scheme. It is an open question to investigate the effect of these characteristics and diversities on the two categories of hybrid scheme.

Table XII: Comparison between Replication and Erasure Coding schemes.

Schemes	Availability	Durability	Recovery	Storage Overhead	Repair Traffic	Read/Write latency
Replication	High	Low	Easy	>1X	=1X	Low for hot-spot objects with small size
Erasure Coding	Low	High	Hard	<1X	>1X	Low for cold-spot objects with large size

Table XIII: Comparison between the state-of-the-art projects using different redundancy schemes.

Projects	Redundancy Scheme	Contributing Factors	Objective(s)
DepSky [Bessani et al. 2011]	Replication	AR†, Pr	High availability, integrity, and confidentiality
Spanner [Wu et al. 2013]	Replication	OS, AR, Pr	Cost optimization and guaranteed availability
CosTLO [Wu et al. 2015a]	Replication	OS, AR, Pe	Optimization of variance latency
SafeStore [Kotla et al. 2007]	Erasure coding	AR, Pr	Cost optimization
RACS [Abu-Libdeh et al. 2010]	Erasure Coding	AR	Cost optimization and vendor-lock in
HAIL [Bowers et al. 2009]	RAID technique	n/a	High availability and integrity
NCCloud [Chen et al. 2014]	Network Codes	n/a	Recovery cost optimization of lost data
CDStore [Li et al. 2015]	Reed-Solomon	n/a	Cost optimization, security and reliability
CAROM [Ma et al. 2013]	Hybrid (RAC)	AR	Cost optimization
CHARM [Zhang et al. 2015]	Hybrid (ROC)	AR, Pr	Cost optimization and guaranteed availability
HyRD [Mao et al. 2016]	Hybrid (ROC)	OS, Pr, Pe	Cost and latency optimization

† n/a: (not applicable), OS: (object size), AR: (access rate), Pr: (price), and Pe (performance).

## 5. DATA CONSISTENCY

Data consistency means that data values remain the same for all replicas of a data item after an update operation. It is investigated in three main aspects, *level*, *metric*, and *model*, as shown in Fig. 8. This section first describes different consistency levels and their pros and cons (§5.1). Then it defines the consistency metrics to determine how much a consistency semantic/model is stronger than another (§5.2). Finally, it discusses consistency models from the user-perspective (§5.3) and from the data store perspectives— weak consistency (§5.4, §5.5) and adaptive consistency (§5.6.1).

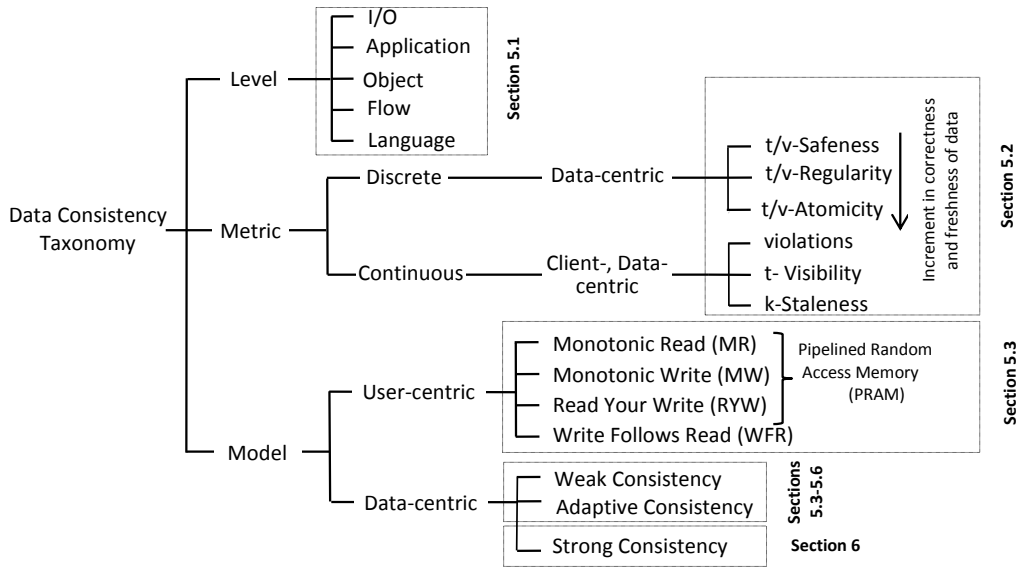


Fig. 8: Data consistency taxonomy

### 5.1. Consistency level

Distributed data stores rely on different levels of data consistency [Alvaro et al. 2013], as shown in Fig. 8.

*I/O-level* consistency allows a clear separation between low-level storage and application logic. It simplifies the development of the application and the complexity of distributed programming. However, I/O-level consistency requires conservative assumptions like concurrent write-write and read-write on the same data, resulting in inefficiency. It should also execute writes and reads in a serial order due to its unawareness of the application semantics. *We focus on this level of consistency in this paper.*

*application-level* consistency exploits the semantics of the application to ensure the concreteness of invariants<sup>7</sup> without incurring the cost of coordination among operations. Thus, it imposes a burden on the developers and sacrifices the generality/reusability of the application code.

*Object-level* consistency makes a trade-off between *efficiency* and *reusability* respectively degraded by I/O- and application-level consistency. It provides the convergence of replicas to the same value without any need of synchronization across replicas via Conflicted-free Replicated Data Types (CRDTs) [Shapiro et al. 2011] in which the value of objects can change in an *associative*, *commutative*, and *idempotent* fashion. Though object-level consistency removes concerns of the reusability of application-level consistency, it requires mapping the properties of the application to invariants over objects by developers.

*Flow-level* consistency is an extension of object-level consistency and requires a model to obtain both the semantic properties of *dataflow component* and the dependency between interacting components.<sup>8</sup> Some components are insensitive to message order delivery as a semantic property,<sup>9</sup> and they are *confluent* and produce the same set of outputs under all ordering of their inputs [Alvaro et al. 2014]. Flow-level consistency demands manual definition for confluent components, resulting in error-prone. But it is more powerful than object-level consistency and it has more applicability than object- and language-level consistency. Indigo [Balegas et al. 2015] exploits flow-level consistency based on confluent components which contain the application-specific correctness rules that should be met.

Finally, *language-level* consistency integrates the semantics and dependencies of the application and maintains a long history of invariants to avoid distributed coordination across replicas. The CALM principle [Alvaro et al. 2011] shows a strong connection between the need of distributed coordination across replicas and *logical monotonicity*. Bloom language [Alvaro et al. 2011] deploys this principle and translates *logical monotonicity* into a practical program that is expressible via *selection*, *projection*, *join*, and *recursion* operations. This class of program, called *monotonic program*,<sup>10</sup> provides output as it receives

<sup>7</sup>The term invariant refers to a property that is never violated (e.g., primary key, foreign key, a defined constraint for the application-e.g., an account balance  $x \geq 0$ ).

<sup>8</sup>A component is a logical unit of computing and storage and receives streams of inputs and produces streams of outputs. The output of a component is the input for other components, and these streams of inputs and outputs implement the flow of data between different services in an application.

<sup>9</sup>The semantic property is defined by application developers. For example, developers determine confluent and non-confluent path between components based on analysis of a component's input/output behavior.

<sup>10</sup>Non-monotonic program contains aggregation and negation queries, and this type of program is implementable via block algorithms that generate output when they receive the entire inputs set.



input elements, and thus guarantees eventual consistency under any order of inputs set. Unlike Bloom, QUELEA language [Sivaramakrishnan et al. 2015] maps operations to a fine-grained consistency levels such as *eventual*, *causal*, and *ordering* and transaction isolation levels like *read committed* (RC), *repeatable read* (RR) [Berenson et al. 1995], and *monotonic atomic view* (MAV) [Bailis et al. 2013].

## 5.2. Consistency metric

This determines how much a consistency model is stronger than another and is categorized into *discrete* and *continuous* from data store perspective (i.e., data-centric) and user perspective (i.e., client-centric).

*Discrete metrics* are measured with the maximum number of time unit (termed by *t-metric*) and data version (termed by *v-metric*). As shown in Fig. 8, they are classified into three metrics [Anderson et al. 2010]: (i) *Safeness* mandates that if a read is not concurrent with any writes, then the most recent written value is retrieved. Otherwise, the read returns any value. (ii) *Regularity* enforces that a read concurrent with some writes returns either the value of the most recent write or concurrent write. It also holds *safeness* property. (iii) *Atomicity* ensures the value of the most recent write for every concurrent or non-concurrent read with write.

*Continuous metrics*, shown in Fig. 8, are defined based on *staleness* and *ordering*. The former metric is expressed in either *t-visibility* or *k-staleness* with the unit of probability distribution of *time* and *version lag* respectively. The latter one is measured as (i) *the number of violations* per time unit from data-centric perspective and (ii) *the probability distribution of violations* in the forms of MR-, MW-, RYW-, WFR-violation from client-centric perspective as discussed later.

## 5.3. Consistency model

This is classified into two categories: *user-* and *data-centric* which respectively are vital to application and system developers [Tanenbaum and Steen 2006].

*User-centric consistency* model is classified into four categories as shown in Fig. 8. *Monotonic Read* (MR) guarantees that a user observing a particular value for the object will never read any previous value of that object afterwards [Yu and Vahdat 2002]. *Monotonic Write* (MW) enforces that updates issued by a user are performed on the data based on the arrival time of updates to the data store. *Read Your Write* (RYW) mandates that the effects of all writes issued by users are visible to their subsequent readers. *Write Follows Read* (WFR) guarantees that whenever users have recently read the updated data with version  $n$ , then the following updates are applied only on replicas with a version  $\geq n$ . Pipelined Random Access Memory (PRAM) is the combination of MR-, MW-, and RYW-consistency and guarantees the serialization in both of reads and writes within a session. Brantner et al. [2008] designed a framework to provide these client-centric consistency models and the atomic transaction on Amazon S3. Also, Bermbach et al. [2011] proposed a middle-ware on eventually consistent data stores to provide MR- and RYW-consistency.

*Data-centric consistency* aims at coordinating all replicas from the internal state of data store perspective. It is classified into three models. *Weak consistency* offers low latency and high availability in the presence of network partitions and guarantees *safeness* and *regularity*. But it causes a complicated burden on the application developers and caters the user with the updated data with a delay time called *inconsistency window* (ICW). In contrast, *strong consistency* guarantees simple semantics for the developer and *atomicity*. But it suffers from long latency which is eight times more than that of weak consistency, and consequently its performance in reads and writes diverges by more than two orders of magnitude [Terry et al. 2013]. *Adaptive consistency* is switching between a range of weak and strong consistency models based on the application requirements/constraints like availability, partition tolerance, and consistency.

The consistency (C) model has a determining effect on achieving availability (A) and partition tolerance (P). Based on the CAP theorem [Gilbert and Lynch 2002], data stores provide only two of these three properties. In fact, data stores offer only CA, CP, or AP properties, where CA in CAP is a better choice within a data store due to rare network partition, and AP in CAP is a preferred choice across data stores. Recently, Abadi [Abadi 2012] redefined CAP as PACELC in order to include latency (L) that has a direct influence on monetary profit and response time, especially across data stores, where the latency between DCs might be high. The term PACELC means that if there is a network partition (P) then there is a choice between A and C for designers, else (E) the choice is between L and C. Systems like Dynamo and Riak leave strong consistency to achieve high availability and low latency. Thus, they are PA/EL systems. Systems with full ACID properties attain stronger consistency at the cost of lower availability and higher latency. Hence, they are PC/EC systems. See Table I in the Appendix for more examples.

## 5.4. Eventual consistency

In this section, we first define the eventual consistency model (§5.4.1). Then, we discuss how this model is implemented and describe how the conflicts that arise from this model are solved (§5.4.2).

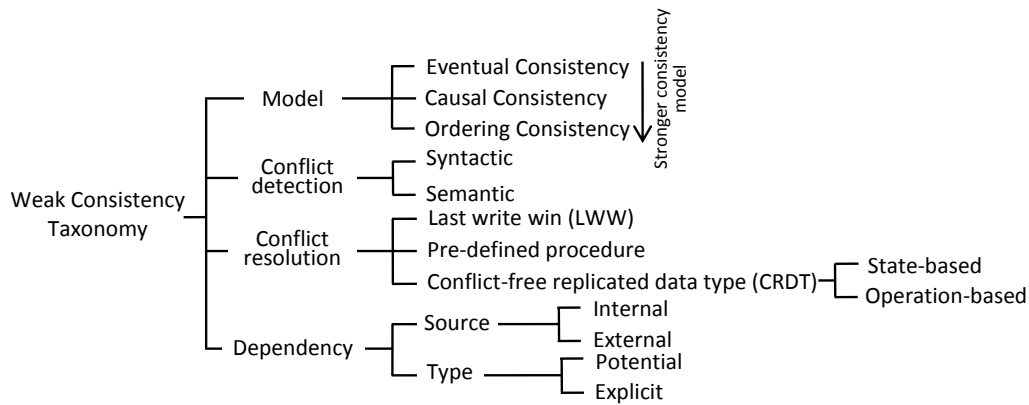


Fig. 9: Weak consistency taxonomy

**5.4.1. Definition.** *Eventual consistency* is defined as all replicas eventually converge to the last update value. It purely supports *liveness* which enforces that all replicas eventually converge based on the operations order, while lacking *safety* which determines the correct effects of operations and leads to incorrect intermediate results. The safety property is assessed in terms *t-visibility* and *k-staleness* as *inconsistency window* (ICW) which is affected by the communication latency, system load, and replicas number. Bailis et al. [2012] analytically predicted the value of ICW via probabilistically bounded staleness (PSB) based on the quorum settings (see Section 4.1.5) for quorum-based data stores. Wada et al. [2011] also experimentally measured ICW with different configurations in terms of readers and writers in different threads in the same/different processes on the same/different VMs for S3, Google App Engine (GAE), and SimpleDB. They concluded that S3 and GAE are “good enough” in data freshness with eventual consistency option, while SimpleDB is not. Zhao et al. [2015] also experimentally measured ICW on the virtualization-based DaaS and concluded that the workload characteristics has more effects on ICW as compared to communication latency across geographical locations. It is worth to note that neither NoSQL databases nor DaaS guarantee a specific SLA in terms of data freshness.

**5.4.2. Implementation.** Eventual consistency-based data stores employ optimistic/lazy replication in which (i) the operation is typically submitted to the closest replica and logged/remembered for the propagation to other replicas later, and (ii) replicas exchange the operation or the effect of operation among each other via *epidemic/gossip protocol* [Demers et al. 1987] in the background [Saito and Shapiro 2005]. Operations are partially ordered by deploying vector clocks. This leads to data *conflicts* which happen as operations are simultaneously submitted to the same data in multi-master systems.

There are four approaches to deal with conflicts. (i) *Conflict detection* approaches strengthen the application semantic and avoid the problems arising from ignoring conflicts. These approaches are classified into *syntactic* and *semantic* [Saito and Shapiro 2005]. The syntactic approach relies on logical or physical clock, whereas the semantic approach works based on the semantic knowledge of operations such as invariants, commuting updates (i.e., CRDTs), and pre-defined procedure. (ii) *Conflicts prohibition* is attainable via blocking or aborting operations and using a single master replica, which comes at the expense of low availability. (iii) *Conflicts ignorance* and *reduction* are achievable by the following conflict resolution techniques (Fig. 9) to guarantee *safety*.

(1) *Last Write Win (LWW)* [Thomas 1979] ignores conflicts, and the update with the highest timestamp is accepted (e.g., Riak (by default), SimpleDB, S3, and Azure Table<sup>11</sup>). It causes lost updates (i.e., updates with less timestamp as compared with winner update) and the violation of expected semantics.

(2) *Pre-defined procedure* merges two versions of a data item to a new one according to application-specific semantic as used in Dynamo. The merges must be *associative* and *commutative* for guaranteeing eventual/causal consistency. Albeit the pre-defined procedure solves conflicts without the need of the total order, it is error-prone and lacks generality. Some data stores use *application-specific precondition* (i.e., a condition or predicate that must always be true just prior to the execution of other conditions) to determine *happened-before dependencies* among requests when the causal consistency model comes as a need.

(3) *Conflict-free replicated data types (CRDTs)* avoid the shortcomings of the above approaches and provide eventual consistency in the presence of node failure and network partition without cross replicas coordination. CRDTs enforce the convergence to the same value after all updates are executed on replicas, and they are either *operation-based* or *state-based* [Shapiro et al. 2011].

<sup>11</sup><https://azure.microsoft.com/en-us/services/storage/>

In *state-based* CRDT, the local replica is first updated and then the modified data is transmitted across replicas. State-based CRDT pursues a partial order  $\leq_v$  (e.g., integer order) with *least upper bound* (LUB)  $\sqcup_v$  (e.g., maximum or minimum operation between integer numbers) that guarantees *associative* (i.e.,  $(a_1 \sqcup_v a_2) \sqcup_v a_3 = a_1 \sqcup_v (a_2 \sqcup_v a_3)$ ), *commutative* (i.e.,  $a_1 \sqcup_v a_2 = a_2 \sqcup_v a_1$ ), and *idempotent* (i.e.,  $a_1 \sqcup_v a_1 = a_1$ ) properties, for each value of object  $a_1, a_2, a_3$  (e.g., integer numbers). Such CRDT is called *Convergent Replicated Data Type* (CvRDT) and is used in Dynamo and Riak. CvRDT can tolerate out-of order, repeatable, and lost messages as long as replicas reach the same value. Thus, CvRDT achieves eventual consistency without any coordination across replicas, but it comes at the expense of monetary cost and communication bottleneck for transferring large objects particularly across data stores. Almeida et al. [2015] addressed this issue by propagating the effect of recent update operations on replicas instead of the whole state; meanwhile all properties of CvRDT are maintained. As an example of CvRDT, consider Grow-only set (G-set) that supports only *union* operations. Assume a partial order  $\leq_v$  on two replicas of G-set  $S_1$  and  $S_2$  is defined as  $S_1 \leq_v S_2 \iff S_1 \subseteq S_2$  and *union* is performed as  $S_1 \cup S_2$ . Since the *union* operation preserves the mentioned three properties, G-set is a CvRDT.

In *operation-based* CRDT, first the update is applied to the local replica, and then is asynchronously propagated to the other replicas. Operation-based CRDT demands a reliable communication network to submit all updates to every replica in a delivery order  $\leq_v$  (specified by data type) with *commutative* property [Shapiro et al. 2011], as utilized in Cassandra. If all concurrent operations are commutative, then any order of operations execution converges to an identical value. Such data type is called *Commutative Replicated Data Type* (CmRDT) and is more useful than CvRDT in terms of data transferring for applications that span write-intensive replicas across data stores. This is because CmRDT demands less bandwidth to transfer *operation* across replicas, as compared to CvRDT that transfers the effect of operation. For instance, G-set is also CmRDT because *union* is commutative. Similar to CvRDT, CmRDT allows the execution of updates anywhere, anytime, and any order, but they have a key shortcoming in guaranteeing integrity constraints and invariants across replicas.

## 5.5. Causal and Causal+ consistency

We first introduce a formal definition of causal and causal+ consistency models (§5.5.1), followed by a description of the source and type of dependencies found in this model (§5.5.2). We then discuss the state-of-the-art projects supporting these models of consistency (§5.5.3).

**5.5.1. Causal consistency definition.** Causal consistency maintains the merits of eventual consistency, while respecting to the causality order among requests applied to replicas. It is stronger and more expensive than eventual consistency due to tracking and checking dependencies. It defines Lamport's "*happens-before*" relation [Gray and Lamport 2006] between operations  $o_1$  and  $o_2$  as  $o_1 \rightsquigarrow o_2$ . Potential causality  $o_1 \rightsquigarrow o_2$  maintains the following rules [Ahmad et al. 1995]. *Execution thread*: If  $o_1$  and  $o_2$  are two operations in the same thread of execution, then  $o_1 \rightsquigarrow o_2$  if  $o_1$  happens before  $o_2$ . *Read from*: If  $o_1$  is a write, and  $o_2$  is a read and returns the value written by  $o_1$ , then  $o_1 \rightsquigarrow o_2$ . *Transitivity*: if  $o_1 \rightsquigarrow o_2$  and  $o_2 \rightsquigarrow o_3$ , then  $o_1 \rightsquigarrow o_3$ . Causal consistency does not support concurrent operations (i.e.,  $a \not\rightsquigarrow b$  and  $b \not\rightsquigarrow a$ ). According to this definition, the write operation happens if all write operations having causal dependency with the given write have occurred before. In other words, if  $o_1 \rightsquigarrow o_2$ , then  $o_1$  must be written before  $o_2$ . Causal+ is the combination of *causal* and *convergent conflict resolution* to ensure *liveness* property. This consistency model allows users locally receive the response of read operations without accessing remote data store, meanwhile the application semantics are preserved due to enforcing causality on operations. However, it degrades scalability across data stores for write operations because each DC should check whether the dependencies of these operations have been satisfied or not before their local commitment. This introduces a trade-off between *throughput* and *visibility*. *Visibility* is the amount of time that a DC should wait for checking the required dependencies among the write operations in the remote DC, and can be influenced by *network latency* and *DC capacity for checking dependencies*.

**5.5.2. Dependency source and type.** Dependencies between operations are represented by a graph in which each vertex represents an operation on variables and each edge shows the causality of a dependency between two operations. The source of dependencies can be *internal* or *external* [Du et al. 2014a]. The former refers to causal dependencies between each update and previous updates in the same session, while the latter relates to causal dependency between each update and updates created by other sessions whose values are read in the same session. COPS [Lloyd et al. 2011], Eiger [Lloyd et al. 2013], and Orbe [Du et al. 2013] track both dependency sources. Dependency types can be either *potential* or *explicit* for an operation (as in Eiger and ChainReaction [Almeida et al. 2013]) or for a value (as in COPS and Orbe). Potential dependencies capture all possible influences between data dependencies, while explicit dependencies represent the semantic causality of the application level between operations. The implementation of potential dependencies in modern applications (e.g., social networks) can produce large metadata in

size and impede scalability due to generating large dependencies graph in the degree and depth. The deployment of explicit dependencies, as used in Indigo [Balegas et al. 2015], alleviates these drawbacks to some extent, but it is an ad-hoc approach and cannot achieve the desired scalability in some cases (e.g., in social applications). This deployment is made more effective with the help of garbage collection, as used in COPS and Eiger, in which the committed dependencies are eliminated and only the nearest dependencies for each operation are maintained.

**5.5.3. Causally consistent data stores.** Causal consistency recently received significant attention in the context of Geo-replicated data stores. **COPS** [Lloyd et al. 2011] provides causal+ consistency by maintaining metadata of causal dependencies. In COPS, a read is locally submitted, and the update operations become visible in a DC when their dependencies are satisfied. COPS supports a causally consistent read-only (RO) transactions which return the version of objects. To do so, it maintains full list of all dependencies and piggyback them when a client issues read operations, as opposed to maintaining the nearest dependencies for providing causal+ consistency. **Eiger** [Lloyd et al. 2013] supports the same consistency model for the column-family data model and provides RO and WO transactions. It maintains fewer dependencies and eliminates the need for garbage collection as compared to COPS. **ChainReaction** [Almeida et al. 2013] uses a variant of chain replication [van Renesse and Schneider 2004] to support causal+ consistency. Contrary to COPS requiring each DC to support *serializability* (see §6), it uses two logic clocks (LCs) for each object: the first one is a global LC, and the second one is local LC that determines which local replica can provide causal consistency. ChainReaction provides RO transactions the same as COPS whilst averting 2 round trip times (RTTs) (as required in COPS) by deploying a sequencer in each DC. Although the sequencer reduces the number of RTT (at most 1 RTT) in the case of RO transactions, it reduces scalability and increases the latency for all updates by 1 RTT within DC. **Orbe** [Du et al. 2013] deploys DM-protocol (exploiting a matrix clock) to support basic reads and writes and DM-clock protocol (deploying physical clock to track dependencies) to provide RO transactions the same as those in ChainReaction that completes in 1 RTT. Orbe avoids the downsides of ChainReaction due to not using a centralized sequencer. **GentleRain** [Du et al. 2014b] eliminates dependencies checking and reduces the metadata piggybacked to each update propagation. This innovation achieves throughput analogous with that in eventual consistency and reduces storage and communication overheads. To achieve such aims, GentleRain uses physical time and allows a DC to make the update visible if all partitions within the DC have seen all updates up to the remote update timestamp. Nevertheless, this technique deteriorates updates visibility in remote DCs. **Bolt-on causal consistency** [Bailis et al. 2013] inserts a layer between clients and data stores to provide causal consistency according to the semantics of the application, not the deployment of LCs or physical clocks. The discussed projects are summarized in Table III-Appendix.

## 5.6. Ordering and Adaptive-level consistency

We first define ordering consistency model and how it is provided. We then discuss projects that enable application providers to switch between a range of consistency models based on their requirements.

**5.6.1. Ordering consistency.** As discussed earlier, eventual consistency applies the updates in different orders at different replicas and causal consistency enforces partial ordering across replicas. In contrast, *ordering consistency* –also called *sequential consistency*– provides a global ordering of the updates submitted to replicas by using a logical clock to guarantee monotonic reads and writes. In fact, *ordering consistency* mandates a read operation from a location to return the value of the last write operation to that location. Ordering consistency is guaranteed through deploying *chain replication* [van Renesse and Schneider 2004] or a master replica which is responsible for ordering writes to an object and then propagating the updates to slave replicas as provided per key in PNUTS.

**5.6.2. Adaptive-level consistency.** *Adaptive-level consistency* switches between weak and strong consistency models based on the requirements of applications to reduce response time and monetary cost. The following data stores leverage adaptive-level consistency.

**CRAQ** [Terrace and Freedman 2009] switches between strong, causal, and eventual consistency for reading objects replicated in a chain replication topology. CRAQ's current implementation relies on placing chains within a DC with the capability of stretching on multiple DCs. It provides a single row (object) consistency and *mini-transactions* that update multiple objects in a single or multiple chains. **Pileus** [Terry et al. 2013] offers a broad spectrum of consistency levels between strong and eventual based on SLAs like latency. It supports all kinds of transactions, but does not scale well because all writes are assigned to the primary replica without automatic movement to other DCs in the face of workload changes. **Gemini** [Li et al. 2012] switches between strong and eventual consistency based on blue and red operations which are respectively executed on different and same order at different DCs. It also introduces the concept of shadow operations to increase blue operations for improving response time and throughput. The consistency level of such operations can be manually assigned by developers based on a specified

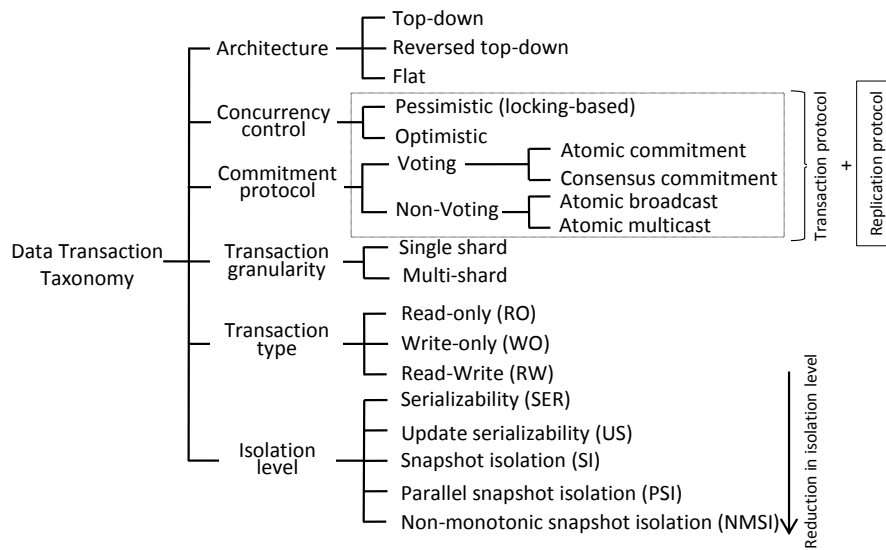


Fig. 10: Data transaction taxonomy

method as in Gemini or based on the SLA as in Pileus. These methods are non-trivial and cumbersome, and **SIEVE** [Li et al. 2014] alleviates this by having an automatic assignment mechanism that exploits application code, invariants, and CRDTs. SIEVE incurs low run-time overheads, but it lacks scalability as applications code becomes large. The discussed projects are summarized in Table IV-Appendix.

## 6. TRANSACTIONS

A *transaction* is defined as a set of reads and writes followed by a *commit* if the transaction is completed or an *abort* otherwise. A transaction has four properties: (i) *atomicity* guarantees *all-or-nothing* outcomes for the set of operations, (ii) *consistency* ensures any transaction transits the database from one valid state to another, (iii) *isolation* ensures the concurrent execution of transactions leads to a database state that would be obtained if transactions were executed one after the other, (iv) *durability* means that once a transaction is committed, all the changes have been recorded to a durable data store. These four properties are called ACID properties. Of these properties, *isolation* concurrency level (*isolation level* for short) refers to a control mechanism for executing two concurrent transactions accessing the same data item. Isolation level expresses the consistency semantic provided by the transaction. Data stores initially provided eventual consistency which is suitable for certain applications (e.g., social networks), while some classes of applications (e.g., e-commerce) demand strong consistency. Data stores, therefore, shifted to guarantee single row transactions (e.g., SimpleDB [Wada et al. 2011] and PNUTS [Cooper et al. 2008]), single shard transactions (e.g., SQL Azure [Campbell et al. 2010]), and multi-shard transactions where shards are geographically replicated (e.g., Spanner [Corbett et al. 2013]).

*Serializability* as the strongest isolation level is not scalable across data stores because it requires strict coordination between replicas probably located far from each other across data stores. Such isolation level has overheads like latency increment, throughput reduction, and unavailability in the case of network partition. Thus, data stores leave serializability, and they resort to (i) offer weaker isolation levels, (ii) partition and replicate data accessed by transactions at a limited number of servers/DCs, and (iii) exploit techniques to optimize coordination as the key requirement to maximize scalability, availability, and performance. Ardekani et al. [2013] stated the following criteria for achieving these purposes across data stores. (1) *Wait-free query* (WFQ) means that a read-only (RO) transaction always commits without synchronizing with other transactions. (2) *Genuine Partial Replication* (GPR) decreases the synchronization and computational time. (3) *Minimal commitment synchronization* is achievable if synchronization is not avoidable (e.g., two transactions with write-write conflicts must be synchronized). (4) *Forward freshness* is attainable if a transaction is allowed to read the committed object version after the transaction starts. Unfortunately, classical concurrency protocols (e.g., 2PL+2PC) are not scalable and purely applicable across the data stores to achieve the above goals. To understand how to effectively deploy classical protocols, we study the main aspects of transactional data stores in the following section.

### 6.1. Transactional data stores

This section details a taxonomy of transactional data stores as shown in Fig. 10.



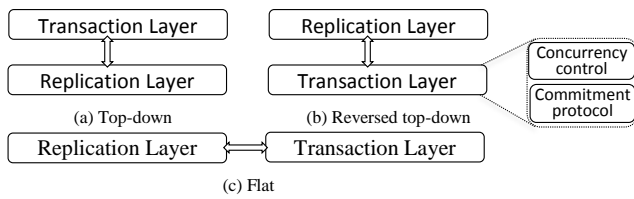


Fig. 11: Transaction architectures in cloud storage

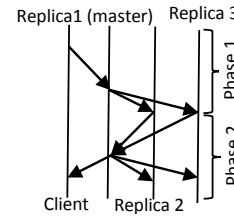


Fig. 12: Execution of basic Paxos

**6.1.1. Architecture.** To build an ACID transactional data store, a stacked layer consisting of *transaction* and *replication* layers is used. The transaction layer consists of, shown in Fig. 11b), (i) *concurrency control mechanism* [Kung and Robinson 1981] that schedules a transaction within each shard to guarantee *isolation*, and (ii) an *atomic commitment protocol* that coordinates distributed transactions across shards to provide *atomicity*. The replication layer is responsible for synchronizing replicas in the case of strong consistency. Fig. 11 shows different architectures that indicate how transaction and replication layers are configured [Agrawal et al. 2013]. Notable systems like Spanner [Corbett et al. 2013] preserves the *top-down* architecture (Fig. 11a), while Replicated Commit [Mahmoud et al. 2013] follows the *reversed top-down* architecture (Fig. 11b) to achieve lower latency. In the *flat* architecture (Fig. 11c), both layers access the storage layer, and there is no clear border between them as deployed in Megastore [Baker et al. 2011]. The separation of the two layers in the architecture brings advantages like modularity, clarity of semantic, and less message-exchange.

**6.1.2. Concurrency control.** This schedules transactions accessing the same data so that the *isolation* property is guaranteed and is categorized into *pessimistic (locking-based)* and *optimistic*. In the *pessimistic* protocol, a transaction first locks (e.g., via 2 Phase Lock (2PL)) the shared data accessed by concurrent transactions, and then operates on data. This protocol serializes transactions upon occurrence of conflict incidences and requires a deadlock detection and resolution mechanism. Therefore, it increases response time especially across data stores, reduces the throughput of data store (expressed as the number of committed transactions per time unit), hurts availability if the lock is held by a failed node, and suffers from *thrashing* (i.e., when the number of transactions is high, many of them become blocked and only few are in progress). Nevertheless, it eases write-write conflicts detection without maintaining the transaction meta-data and works best when conflicts among transactions are short running. On the contrary, in *optimistic* protocol (optimistic concurrency control (OCC)), all transactions are executed separately though each of them before commitment needs to pass a *certification procedure* in which the write set of transaction  $T$  that is being validated against the read and write sets of other active transactions in the system. Thus, any read-write or write-write conflicts result in OCC to be aborted. To avoid such procedure, MVOCC [Vo et al. 2012] combines MVCC and OCC to detect conflicts via the version number of data. Thus, MVOCC always commits RO transactions, while uses the version number of data to check their conflicts for update transactions. However, OCC provides low latency at the cost of weak concurrency control, fits when transactions are long running and avoids drawbacks arising from the pessimistic protocol though it wastes resources due to restarting transaction and requiring exclusive lock during final 2 Phase Commit (2PC) [Agrawal et al. 1987]. MaaT [Mahmoud et al. 2014], as a re-design of OCC, eliminates locks during 2PC to reduce aborts rate, avoids multi-version concurrency control to make efficient use of memory, and improves throughput; meanwhile, MaaT maintains the *no-thrashing* property of OCC.

**6.1.3. Commitment protocol.** This refers the way transactions terminate while guaranteeing *atomicity* and *consistency*. It is generally categorized into two types as below:

(1) *Voting* protocols agree on updates based on voting all replicas as in *atomic commitment protocol* (e.g., 2PC) or quorum of replicas as in *consensus commitment protocol* (e.g., Paxos commitment [Gray and Lamport 2006]). In 2PC protocol, a node is designated as the *coordinator/master* node and other nodes in the system as *participants*, *cohorts*, or *workers*. This protocol consists of (i) a *voting phase* in which a coordinator node prepares all participant nodes to take either commit (“yes”) or abort (“no”) vote on the transaction, and (ii) *commit phase* in which the coordinator decides whether to commit or abort the transaction based on vote of participant nodes in the previous phase, and then notifies about its decision to all the participant nodes. This protocol suffers from blocking and non-resilience to node failures. By contrast, Paxos protocol is fault-tolerant, non-blocking, and a consensus algorithm for achieving on a single value among a set of replicas and tolerates a variety of failures like duplicated, lost, and reordered messages as well as the failure of nodes. In the deployed applications across data stores, storage and computing nodes take the responsibility of four basic roles in Paxos: *client* are application servers, *proposers* are coordinators, *acceptors* are storage nodes, and all nodes are *learners*. As shown in Fig. 12, this protocol

requires 2 RTTs to reach consensus on a value and consists of two phases: the first establishes the coordinator/master for an update for specific record issued by a client, and the second tries to reach a consensus value across a majority of acceptors and writes the specified value for a specific record.

(2) *Non-voting* protocols work based on a group communication (GC) primitive in which, unlike 2PC, all nodes receive the certificate vote and then locally perform a commit or abort operation. The GC is mainly categorized into two primitives. *Atomic broadcast* [Défago et al. 2004] propagates messages to *all* DCs and ensures that all DCs agree on the set of received messages and their order. It offers serializability in the case of full replication but hampers scalability, which can be relieved via pairing it with genuine partial replication (GPR). By contrast, *atomic multicast* [Schiper et al. 2009] propagates the message to a subset of DCs which can be *genuine* or *non-genuine* [Guerraoui and Schiper 2001]. Genuine protocols are expensive in terms of delivering message and receiving acknowledgement, while non-genuine ones can deliver messages in 1 RTT and can work well except when the load and the number of DCs are high in the system [Schiper et al. 2009]. Moreover, Genuine and non-genuine protocols are identical in *minimality property* which implies messages are only sent to nodes— servers or DCs — hosting the desired replicas. Note that, atomic multicast implementation using atomic broadcast protocol does not satisfy the minimality property, and thus it is *non-genuine*.

**6.1.4. Replication protocol.** Data stores rely on a *replication protocol* to guarantee a serial ordering of write commits of different transactions. They use replication protocol like basic Paxos [Lamport 1998], View-stamped [Oki and Liskov 1988] (equivalent to Multi-Paxos-based) or atomic broadcast [Défago et al. 2004]. These protocols are expensive at (i) throughput as a function of the load on a bottleneck replica (e.g., leader replica), and (ii) latency as a function of the number of message delays in the protocol— i.e., from when client sends a request to replica until it receives a reply. As shown in Fig. 12, the message flow in a leader-based Paxos is: client→leader→replicas →leader→client; thus the latency (i.e., *message delays*) is 4 and throughput (i.e., *bottleneck at message— here master replica is a bottleneck*) is  $2n$ , where  $n$  is the number of replicas. Fast Paxos [Lamport 2006] reduces this latency by sending requests from the client to replicas instead of through a distinguished replica as a leader. It reduces latency to 3 (client→replicas →leader→client) at the cost of requiring larger quorum sizes as compared to that for Paxos. Generalized Paxos [Lamport 2005] is an extension of Fast Paxos and exploits *commutative property* between operations to commit a request in two message delays. EPaxos [Moraru et al. 2013] is another replication protocol which is built on Fast Paxos and Generalized Paxos to achieve low latency and high throughput. Unlike these protocols, Inconsistent replication (IR) [Zhang et al. 2015] is fault-tolerant without guaranteeing any consistency.

As discussed, the fundamental difference between replication protocols is an extra RTT or a large quorum to order conflicts. Replication protocols recently provide *ordering* in network layer and leave reliability to the replication layer. This results in a fewer message delay and higher throughput. Speculative Paxos [Ports et al. 2015] assumes a best effort ordering at the cost of application-level roll-back, while NOPaxos [Li et al. 2016] guarantee ordering in network layer to avoid roll-back/transaction abort. Table XIV summarizes consensus replication protocols in performance metrics.

Table XIV: Comparison between different consensus replication protocols.

Performance metrics	Basic Paxos†	Fast Paxos	Generalized Paxos	Paxos Batching	Inconsistent Replication	Speculative Paxos	NOPaxs
Latency (message delay)	4	3	2	>4	2 4††	2	2
Message at bottleneck	$2n$	$2n$	$2n$	$2n$	2	2	2
Quorum size	$> n/2$	$> 2n/3$	$> 2n/3$	$> n/2$	$> 2n/3 n/2$	$> 3n/4$	$> n/2$

† Multi-Paxos is an optimization for basic Paxos by reserving the master replica for several Paxos instances; thus it can avoid phase 1 to achieve two message delays for agreement on a consensus value.

†† Inconsistent Replication (IR) uses the fast path (slow path) which requires 2 (4) message delays with a quorum size of  $> 2n/3$  ( $> n/2$ ) *consensus* operations. These operations can execute in any order, while inconsistent operations can execute in different order at each replica and IR can complete them in 2 message delays with a quorum size of  $> n/2$ .

**6.1.5. Transaction granularity and type.** One trend in data stores to achieve scalability is to limit the number and location expansion of partitions accessed by a transaction. We call this limitation as *transaction granularity* and classify it as *single shard/partition* and *multi-shard/partition*. The execution and commitment of single shard and multi-shard transactions respectively involves a shard within a server (e.g., SQL Azure [Campbell et al. 2010]) and more than one shard replicated across several servers in a data store or several data stores (e.g., Spanner). Most today's data stores support multi-shard, and some of them pose limitations like the pre-declared read/write set performed by transactions in Sinfonia [Aguilera et al. 2007] and Calvin [Thomson et al. 2012].

Table XV: Anomalies in different isolation levels

Anomalies	SSER	SER	US	SI	PSI	NMSI	MAV	RA	RC	RUC
Dirty Read	N <sup>1</sup>	N	N	N	N	N	N	N	N	Y
Repeatable Reads	N	N	N	N	N	N	N	N	Y	Y
Read Skew	N	N	N	N	N	N	N	Y	Y	Y
Dirty Writes	N	N	N	N	N	N	N	N	N	N
Lost Updates	N	N	N	N	N	N	Y	Y	Y	N
Write Skew	N	N	N	Y	Y	Y	Y	Y	Y	Y
Causality violation	N	N	N	N	Y	Y	Y	Y	Y	Y
Non-Mon. Snap. <sup>2</sup>	N	N	Y	N	Y	Y	Y	Y	Y	Y
True Time violation	N	Y <sup>1</sup>	Y	Y	Y	Y	Y	Y	Y	Y

[1] N: disallowed anomaly and Y: allowed anomaly. [2] Non-Monotonic Snapshot (long fork).

*Transaction type* is divided into three categories. *Read-only* (RO) transactions conduct read operations on any updated replicas without locking; meanwhile the incoming writes are not blocked. The original distributed RO transactions [Chan and Gray 1985], as deployed in Spanner, always take 2 RTTs, and before starting they wait until all the involved servers/DCs guarantee that all the transactions have been committed. In contrast, in Eiger [Lloyd et al. 2013] the RO transactions take 1 RTT by maintaining more metadata. *Write-only* (WO) transactions only contain write operations and deploy concurrency controls to avoid write-write conflicts. *Read-write* (RW) transactions consist of reads and writes and avoid write-read and write-write conflicts by using concurrency controls and commitment protocols.

**6.1.6. Isolation levels.** This refers to the concurrency between transactions. As the isolation level is weaker, the concurrency level is higher but with more anomalies [Adya 1999; Berenson et al. 1995]. In the following, we discuss different isolation levels from strong to weak as shown in Fig. 10.

- *Serializability* (SER) guarantees that every concurrent execution of the committed transactions is equivalent with the serial execution of transactions (i.e., one after another). It is typically implemented via 2PL for full replication, which impedes scalability and increases response time especially across data stores. So, SER is provided in the case of partial replication or genuine partial replication (GPR) (e.g., P-store [Schiper et al. 2010]), in which all transaction types might require a certification procedure and go through a synchronization phase. If so, then the transaction aborts (e.g., Sinfonia and P-store); otherwise, the transaction is *wait-free query* (WFQ) as in S-DUR [Sciascia et al. 2012]. Note that *strict* SER (SSER, also called *linearizability*) deploys the physical clock to order transactions.
- *Update serializability* (US) [Hansdah and Patnaik 1986] provides guarantees analogous to those in SER for update transactions, but it leads RO transactions to observe non-monotonic snapshots<sup>12</sup> (i.e., two RO transactions may observe different order of the committed update transactions). In fact, US makes the support of observing a snapshot equivalent to some serial execution of the partially ordered history of update transactions. Extended (E) US [Peluso et al. 2012], as a stronger variant of US, holds the same semantic not only for committed update transactions but also for those are executing and may abort later due to either write-write or write-read conflicts. This property guarantees that applications do not act in an unexpected manner because of non-serializable snapshots observation.
- *Snapshot isolation* (SI) never aborts RO transactions and guarantees that a transaction reads the most recent version of data through multi-version concurrency control (MVCC) [Kung and Robinson 1981]. However, it blocks write transactions to avoid write-write conflicts and enhances responsiveness of these transactions at the expense of *write skew* anomaly<sup>13</sup> [Berenson et al. 1995] due to ignoring read-write conflicts. It also implements WFQ but not in GPR scheme [Ardekani et al. 2014]. Due to its simple implementation, it is supported by database vendors (e.g., Microsoft SQL Server and Oracle) and most Geo-replicated data stores.
- *Parallel snapshot isolation* (PSI) [Sovran et al. 2011] is similar to SI and suitable to deploy across data stores. In PSI, the commit order of non-conflicting transactions<sup>14</sup> can vary among replicas since the transactions are causally ordered. In fact, a transaction is propagated to other replicas after all transactions committed before it starts. Thus, PSI enforces that transactions observe the local data which may be stale. Unlike SI, PSI neither supports GPR and nor monotonic snapshots of transactions, which is the main hindrance of SI regard to scalability [Ardekani et al. 2013].
- *Non-monotonic snapshot isolation* (NMSI) [Ardekani et al. 2013] works under GPR scheme and preserves WFQ, *minimal commitment synchronization*, and *forward freshness*. NMSI is more scalable and takes any snapshot as compared to PSI, because PSI neither is implementable under GPR and nor is

<sup>12</sup>A snapshot is a logical copy of data consisting of all committed updates and is created when a transaction is committed.

<sup>13</sup>The readers are referred to papers by [Adya 1999; Berenson et al. 1995] on details of anomalies.

<sup>14</sup>Both write-write and write-read conflicts are avoided in SER, SI, and US, while only write-write conflicts are prevented in PSI and NMSI.

preservable for the forward freshness property. The Jessy protocol [Ardekani et al. 2013] implements NMSI under GPR and uses 2PC and the versioning mechanism for concurrency control.

Compared to the above isolation levels, several weaker isolation levels are also implemented in data stores. *Read Uncommitted* (RUC) (as in PNUTS, Dynamo, and BigTable) totally orders writes to each object, and *Read Committed* (RC) disallows transactions to access uncommitted or intermediate version of data items. *Monotonic Atomic View* (MAV) [Bailis et al. 2013] ensures that if a transaction  $T_j$  reads a version ( $v$ ) of an object that transaction  $T_i$  wrote, then a later read by  $T_j$  returns a value whose version  $v' \geq v$ . *Read Atomic* (RA) [Bailis et al. 2014] enforces all or none of each transaction's updates are visible to others. RA prevents all anomalies avoided by RC with the help of multi-versioning and a small metadata per write operations. Unlike RUC, RC, MAV, it also prevents *fractured read* anomaly which happens if transaction  $T_i$  writes variables  $x_m$  and  $y_n$ , then  $T_j$  reads  $x_m$  and  $y_k$ , where  $k < n$ . Both MAV and RA are useful for the applications requiring secondary indexes, foreign key, and materialized view maintenance.

There are several ways to mutually compare isolation levels. The **first** is to understand what anomalies they allow (Table XV). For example, SSER, SER, and US prevent *write-skew* (*short fork*) anomaly due to checking write-write and read-write conflicts. The **second** is that they can be compared based on scalability, which is classified into *highly available transaction* (HAT) and non-HAT [Bailis et al. 2013]. Anomalies like *lost update* and *write skew*, and semantics such as *concurrent update* and *bounds on data recency*—real time insurance on writes/reads—are impossible to be prevented by HAT systems. Thus, SER, SI, PSI, and RR are not HAT-complaint. In contrast, HAT systems preclude *dirty read and write*, and thus RA, RU, and RC are HAT-complaint. The **third** is they can be analyzed via programmatic tool like G-DUR [Ardekani et al. 2014] that implements the most discussed isolation levels.

## 6.2. Coordination mechanisms in transactional data store

This section discusses different mechanisms that coordinate transactions within and across the state-of-the-art data stores.

**6.2.1. Heavy-weight protocols-based transactions.** Transactional data stores usually partition data into *shards* and then replicate each shard across servers for fault-tolerance and availability. To guarantee transactions with strong consistency, they deploy a *distributed transaction protocol* which implements concurrency control via a write-ahead log (log for short). Each shard is designated to a log which is replicated across data stores. This log is divided into *log positions* which are uniquely numbered in ascending order [Agrawal et al. 2013]. To guarantee the serial ordering of write commits of different transactions into logs, transactional data stores commonly integrate a costly *replication protocol* with strong consistency like Paxos. In the Paxos deployment, participants compete to access a *log position* in the log, and only one of them is allowed to commit and others are required to abort. Thus, these data stores enforce strong consistency in both *transaction protocol* which provides a serial ordering of transactions across shards, and *replication protocol* which ensures a serial ordering of writes/reads across replicas in each shard. This redundancy increases latency (message delays) and degrades throughput especially when heavy-weight protocols like 2PL+2PC and OCC+2PC with Paxos are deployed. To alleviate this issue, the following solutions can be leveraged. The first is to use replication protocols with lower latency and higher throughput as discussed before. The second is to replicate and partition data at a limited number of servers/DCs when the data accessed by multi-shard transactions. This solution mandates workload-aware data partitioning approach rather than the commonly used data partitioning mechanisms like *random* (*hash-based*), *round robin*, and *key range*. This approach is classified into *scheme-dependent*, which is coarse-grained and unsuitable for dynamic workload, and *scheme-independent* [Kamal et al. 2016] that is fine-grained and deploys graph partitioning to cater a range of workloads in particular dynamic workload.

**6.2.2. Contemporary techniques.** Sole reliance on heavy-weight protocols deteriorates response time and causes *spurious coordination* as a transaction unnecessarily delays or reorders the execution of transactions based on their natural arrival order. To relieve these drawbacks, the following contemporary techniques are used. The first is *dependency-tracking technique* and is used in the *linear transactions protocol* [Escriva et al. 2012]. This protocol allows transactions to locally commit on the server. ROCOCO [Mu et al. 2014] uses this technique to reorder interfered transactions rather than aborting or blocking them. Warp [Escriva et al. 2015] uses ROCOCO to improve throughput and latency as compared to 2PL and OCC. The second technique is *transaction decomposition*. Lynx [Zhang et al. 2013b] uses this technique and decomposes each transaction to several chains based on the application semantics and the table scheme by using transaction chopping [Shasha et al. 1995]. Homeostasis protocol [Roy et al. 2015] also dynamically extracts application semantics to determine acceptable upper-bound inconsistency of the system. This helps it to locally execute transactions across DCs without communication as long as inconsistency does not violate the correctness of transactions. The third technique is the following set of *logical properties* which minimizes or exempts coordination across transactions: (i) *Commutativity* guarantees reordering

operations which do not influence the outcome of transactions [Lloyd et al. 2013]. *Monotonic programs* ensure deterministic outcome for any order of operations on objects [Alvaro et al. 2011] and support operations such as *selection*, *projection*, and *join* through a declarative language like Bloom. (iii) *CRDTs*, as already discussed, guarantee convergent outcomes irrespective of the order of updates applied on each object. (iv) *Invariants* like primary keys, foreign keys, and integrity constraints reduce coordination between transactions including *confluent* operations that can be executed without coordination with others. However, contemporary techniques are error-prone, ad-hoc, and some of them, like monotonic programs and CRDTs, are not well developed. Tables V and VI in appendix compare the state-of-the-art projects in several aspects and indicates which concurrency mechanism they used.

## 7. DATA MANAGEMENT COST

*Price* is a new and important feature of data stores as compared to traditional distributed systems like cluster and grid computing. Data stores offer a variety of pricing plans for different storage services with a set of performance metrics. Pricing plans offered by data stores are typically divided in two categories [Naldi and Mastroeni 2013]: *bundling price* (also called *quantity discount*) and *block rate pricing*. The first is observed in most data stores (e.g., Google Drive) and is recognized as a non-linear pricing, where unit price changes with quantity to follow *fixed cost* and *per-unit charge*. The second category divides the range of consumption into sub-ranges and in each sub-range unit price is constant as observed in Amazon. This category is a special form of the *multi-part tariffs* scheme in which the fixed price is zero. Note that the standard form of multi-part tariffs consists of a fixed cost plus multi-ranges of costs with constant cost in each range. One common form of this scheme is *two-part tariffs* that are utilized in data stores with a fixed fee for a long term (currently 1 or 3 years) plus a per-unit charge. This model is known a *reserved pricing model* (e.g., as offered by Amazon RDS and Dynamo) as opposed to an *on-demand pricing model* in which there is no fixed fee and its per-unit charge is more than that in the reserved pricing model. All pricing plans offered by the well-known cloud providers follow *concavity property* that implies as the more resources the application providers buy the cheaper the unit price is. The unit price for storage, network, and VM respectively are often GB/month, GB, and instance per unit time.

A cloud provider offers different services with the same functionality while performance is directly proportional to price. For example, Amazon offers S3 and RRS as online storage services but RRS compromises redundancy for lower cost. Moreover, the price of same resources across cloud providers is different. Thus, given these differences, many cost-based decisions can be made. These decisions will become complicated especially for applications with time-varying workloads and different QoS requirements such as availability, durability, response time, and consistency level. To do so, a joint optimization problem of resources cost and the required QoS should be characterized. Resources cost consists of: (i) *storage cost* calculated based on the duration and size of storage the application provider uses, (ii) *network cost* computed according to the size of data the application provider transfers out (reads) and in (writes) to data stores (typically data transfer into data stores is free), and (iii) *computing cost* calculated according to duration of renting a VM by application providers. In the rest of section, we discuss the cost optimization of data management based on a single QoS or multi-QoS metrics, and cost trade-offs.

### 7.1. Cost optimization based on a single QoS metric

Application providers are interested into the selection of data stores based on a single QoS metric so that the cost is optimized or does not go beyond the budget constraint. This is referred to as a cost optimization problem based on a single QoS metric and is discussed as below.

(1) *Cost-based availability and durability optimization*. Availability and durability are measured in the number of nines and achieved by means of usually triplicate replication in data stores. Chang et al. [2012] proposed an algorithm to replicate objects across data stores so that users obtain the specified availability subject to budget constraint. Mansouri et al. [2013] proposed two dynamic algorithms to select data stores for replicating non-partitioned and partitioned objects, respectively, with the given availability and budget. In respect to durability, PRCR [Li et al. 2012] uses the duplicate scheme to reduce replication cost while achieving the same durability as in triplicate replication. CIR [Li et al. 2011] also dynamically increases the number of replicas based on the demanding reliability with the aim of saving storage cost.

(2) *Cost-based consistency optimization*. While most of the studies explored consistency-performance trade-off (§5 and §6), several other studies focused on lowering cost with adaptive consistency model instead of a particular consistency model. The *consistency rationing* approach [Kraska et al. 2009] divides data into three categories with different consistency levels assigned and dynamically switches between them in run time to reduce the resource and the penalty cost paid for the inconsistency of data which is measured based on the percentages of incorrect operations due to using lower consistency level. Bismar [Chihoub et al. 2013] defines the consistency level on operations rather than data to reduce the cost of the required resources at run time. It demonstrates the direct proportion between the consistency level and



cost. C3 [Fetai and Schuldt 2012] dynamically adjusts the level of consistency for each transaction so that the cost of consistency and inconsistency is minimized.

(3) *Cost-based latency optimization*. User-perceived latency is defined as (i) a constant in the unit of RTT, distance, and network hops, or (ii) latency cost that is jointly optimized with monetary cost of other resources. Latency cost allows the latency metric to be changed from a discrete value to continuous one, thereby achieving an accurate QoS in terms of latency constraint and easily making a trade-off between latency and other monetary costs. OLTP (resp. OLAP) applications satisfy the latency constraint by optimization of data placement (resp. data and task placement). This placement has an essential effect on optimizing latency cost or satisfying it as a constraint as well as in reducing resources cost.

## 7.2. Cost optimization based on multi-QoS metric

Application providers employ Geo-replicated data stores to reduce the cost spent on storage, network, and computing under multi-QoS metrics. In addition, they may incur *data migration cost* as a function of the data size transferring out from data store and its corresponding network cost. Data migration happens due to application requirements, the change of data store parameters (e.g., price), and data access patterns. The last factor is the main trigger for data migration as data transits from *hot-spot* to *cold-spot* status (defined in §2) or the location of users changes as studied in “Nomad” [Tran et al. 2011]. In Nomad, the changes in users’ location are recognized based on simple policies that monitor the location of users when they access an object. Depending on the requirements of the applications, cost elements and QoS metrics are determined and integrated in the classical cost optimization problems as linear/dynamic programming [Cormen et al. 2009], k-center/k-median [Jain and Dubes 1988], ski-rental [Seiden 2000], etc. In the following, these features and requirements are discussed.

(1) For a file system deployment, a key decision is to store a data item either in cache or storage at an appropriate time while guaranteeing access latency. Puttaswamy et al. [2012] leveraged EBS and S3 to optimize the cost of file system and they abstracted the cost optimization via a ski-rental problem. (2) For data-intensive applications spanning across DCs, the key decision is which data stores should be selected so that the incurred cost is optimized while QoS metrics are met. The QoS metrics for each data-intensive application can be different; for example, online social applications suffice causal consistency, while a collaborative document editing web service demands strong consistency. (3) For online social network (OSN), the key factor is replica placement and reads/writes redirection, while “social locality” (i.e., co-locating the user’s data and her friends’ data) making reduction in access latency is guaranteed. In OSN, different policies to optimize cost are leveraged: (i) minimizing the number of slave replicas while guaranteeing social locality for each user [Pujol et al. 2010], (ii) maximizing the number of users whose locality can be preserved with a given number of replicas for each user [Tran et al. 2012], (iii) graph partitioning based on the relations between users in OSN (e.g., cosplay [Jiao et al. 2016]), and (iv) selective replication of data across DCs to reduce the cost of reads and writes [Liu et al. 2013]. (4) The emergence of content cloud platforms (e.g., Amazon Cloud-Front<sup>15</sup> and Azure CDN<sup>16</sup>) help to build a cost-effective cloud-based content delivery network (CDN). In CDN, the main factors contributing into the cost optimization are replicas placement and reads/writes redirection to appropriate replicas [Chen et al. 2012].

## 7.3. Cost trade-offs

Due to several storage classes with different prices and various characteristics of workloads, application providers are facing with several cost trade-offs as below.

*Storage-computation trade-off*. This is important in scientific application workloads in which there is a need for the decision on either storing data or recomputing data based on the size and access patterns. Similar decision happens to the privacy preservation context that requires a trade-off between encryption and decryption of data (i.e., computation cost) and storing data [Zhang et al. 2013a]. The trade-off can be also seen in video-on-demand service<sup>17</sup> in which video transcoding<sup>18</sup> is a computation-intensive operation, and storing a video with a variety of formats is storage-intensive. Incoming workload on the video and the performance requirement of users determine whether the video is transcoded on-demand or stored with different formats.

*Storage-cache trade-off*. Cloud providers offer different tiers of storage with different prices and performance metrics. A tier, like S3, provides low storage cost but charges more for I/O, and another tier, like EBS and Azure drive, provides storage at higher cost but I/O at lower cost [Chiu and Agrawal 2010]. Thus, as an example, if a file system frequently issues reads and writes for an object, it is cost-efficient to save the object in EBS as a cache, or in S3 otherwise. This trade-off can be exploited by data-intensive

<sup>15</sup>“Amazon CloudFront”, <https://aws.amazon.com/cloudfront/>.

<sup>16</sup>“Azure CDN”, <https://azure.microsoft.com/en-us/services/cdn/>.

<sup>17</sup><http://techcrunch.com/2009/04/14/hd-cloud-puts-video-formatting-in-the-cloud/>

<sup>18</sup>Video transcoding is the process of converting a compressed digital video format to another.

applications in which the generated intermediate/pre-computed data can be stored in caches such as EBS or memory attached to VM instances.

*Storage-network trade-off.* Due to significant *differences in storage and network costs* across data stores and *time-varying workload* on an object during its lifetime, acquiring the cheapest network and storage resources at the appropriate time of the object lifetime plays a vital role in the cost optimization. Simply placing objects in a data store with either the cheapest network or storage for their whole lifetime can be inefficient. Thus, storage-network trade-off requires a strategy to determine the placement of objects during their lifetime in which the status of objects change from hot-spot to cold-spot and vice versa. This was studied in a dual cloud-based storage architecture [Mansouri and Buyya 2016] as well as in distributed data stores for a limited number of replicas for each object [Mansouri et al. 2017]. This trade-off also comes as a matter in the recovery cost in erasure coding context, where *regenerating* and *non-MDS* codes are designed for this purpose.

*Reserved-on demand storage trade-off.* Amazon RDS and Dynamo data stores offer *on-demand* and *reserved* database (DB) instances and confront the application providers with the fact that how to combine these two types of instances so that the cost is minimized. Although this trade-off received attention in the context of computing resources [Wang et al. 2013], it is worthwhile to investigate the trade-off in regard to data-intensive applications since (i) the workload of these two is different in characteristics, and (ii) the combination of on-demand DB instances and different classes of reserved DB instances with various reservation periods can be more cost-effective. In respect to these cost trade-offs and discussed cost optimization problems with single/multi QoS metrics, Table VII in Appendix summarizes state-of-the-art studies.

## 8. SUMMARY AND FUTURE DIRECTIONS

This section presents a summary and discusses future directions (see Table XVI) of different aspects of data management in cloud-based data stores (data stores for short).

*Intra- and Inter-data store services.* The deployment of OLTP applications within and across data stores brings benefits and challenges for application providers as summarized in Tables II-V. To better achieve the goals of intra-data store (Table II), there are some interesting tracks to explore as identified in Table XVI. However, storing data within data stores faces application providers with challenges as shown in Table III. A way to tackle these challenges, especially *data transfer bottleneck* and performance unpredictability, is to deploy *workload-aware data placement* while considering the topology of data stores rather than randomly placing data in storage nodes of commercial data stores. Furthermore, such solutions require data delivery within a deadline especially for OLTP applications. In respect to another challenge, *data security*, it would be useful to design a fuzzy framework in which acts, legislation, security, and privacy requirements are considered. A complementary approach to these promising solutions is storing data across data stores owned by different vendors. However, the cross-deployment of data stores raises key challenges (Table V). The first one is how to ease the movement across data stores owned by different vendors. This requires the design of a common data model and standard APIs for different cloud databases to help application providers. The second is how to handle network congestion. The ideal solutions for this challenge are to complete data transfer within a *budget* and *deadline* especially for OLTP applications that demand high response time. Last but not the least, concurrency control is another challenge for which the solutions imply to focus on the following data elements as outlined in Fig. 1.

*Data model.* Relational data model provides ACID properties for applications, whereas it impedes scalability. In contrast, NoSQL data model makes data stores more scalable for data-intensive applications. NewSQL takes benefits of both data models. It also mitigates the latency effects of using remote data because a hierarchical scheme of the related data is made, and then the data is placed in servers/DCs which are at the close distance. Thus, it would be very useful to analyse the workload of the application and build a hierarchical scheme of related data. This helps to reduce cross-coordination and confine transactions to a limited number of servers/DCs, and *eases processing complex queries over entities*. Furthermore, both classes of commercial data stores using these data models support weak SLA for a limited performance criteria like availability and durability. Hence, there is a venue to provide better SLA in respect to response time, auto-scaling, monetary cost (Table XVI).

*Data dispersion.* Though data replication has been widely investigated in the field of databases, it is worthwhile to investigate which model, degree, granularity, and propagation of data replication is deployed for achieving specific performance metrics in the cloud-based data stores (Table XVI). In respect to erasure coding, cloud-based data stores require codes which are more efficient in network rather than storage to save monetary cost as deployed across data stores and reduce network congestion as used within data stores. The pure use of each scheme of data dispersion may not be efficient in terms of performance metrics. A better alternative is to deploy a hybrid scheme based on the characteristics of workloads and data stores which were not fully investigated in state-of-the-art projects (Table XIII).

Table XVI: Future Direction for cloud-based data stores

Data store Aspect	Future Directions
Intra- and Inter-data stores services	<ul style="list-style-type: none"> <li>• Designing algorithms to guarantee the time of data retrieval from storage nodes within a range of time via (i) data replication in several storage nodes owned by different vendors, and (ii) redundant requests against replicas while considering optimizing the monetary cost.</li> <li>• Replicating data across storage nodes instead of random replication to avoid <i>correlated</i> and <i>independent</i> failures while considering consistency costs, joining and leaving storage nodes.</li> <li>• Reducing unpredictability of response time in multi-tenant storage services via replicating data in storage nodes and submitting redundant read/write requests against them while considering the status of data, cold- and hot-spot, I/O and CPU load.</li> <li>• Designing auto-scaling mechanisms in which the range of required database instances should be determined based on the workload and the type of instance, not by users who determine currently in commercial data stores.</li> <li>• Designing scheduling algorithms to complete data delivery within deadline and budget by considering the size and price of bandwidth.</li> <li>• Designing a fuzzy and self-learning framework for ensuring data security while considering managerial solutions, acts, legislation, and privacy over data placement within and across data stores.</li> </ul>
Data Model	<ul style="list-style-type: none"> <li>• Analysing the relationship between entities in applications in order to identify their relationship and then placing the associated data in servers/DCs at the close distance</li> <li>• Guaranteeing SLA in terms of response time in both classes of data store via dynamic allocation of bandwidth to requests, selection of adaptive consistency semantic, adaptive replication of data based on their workload.</li> </ul>
Data Dispersion	<ul style="list-style-type: none"> <li>• Borrowing the parallelism in deferred-update replication and abort-free features in state machine replication to improve the scalability and throughput of transactions.</li> <li>• Designing algorithms to replicate data in full or partial degree based on the number of DCs, globality of transactions, and the size of data.</li> <li>• Defining a monetary cost function including latency as a utility to capture it as a continuous value instead of discrete and the cost of storage services, and then applying this cost function in Quorum-based replication approach.</li> <li>• Constructing codes with small repair bandwidth and minimal number of nodes participating in the repair process of the failed chunk for cold objects, and constructing codes for hot objects with the help of queuing and coding theory to read and to construct the failed chunk from several data chunks in parallel.</li> <li>• Utilizing fully the characteristics of workload and data stores in performance and price to leverage hybrid scheme of data dispersion.</li> </ul>
Data consistency and transaction management	<ul style="list-style-type: none"> <li>• Extension of Conflicted-free Replicated Data Types (CRDTs) in operation types and make the richer relation between them in order to exempt the coordination between data replicas.</li> <li>• Analysing the semantic of the application to provide consistency at the language level.</li> <li>• Designing replication protocols with the capability of <i>ordering</i> request against replicas at the <i>network level</i> with the small quorum size in order to increase throughput— see Table XIV.</li> <li>• Designing an adaptive isolation levels for applications based on the required response time and available budgets. As response time is tighter, weaker isolation level is chosen.</li> <li>• Designing adaptive algorithms with the help of graph partitioning techniques and clustering algorithms to reduce distributed transactions, and in turn, improve response time.</li> </ul>
Data management cost	<ul style="list-style-type: none"> <li>• Designing algorithms to make a trade-off between performance criteria like availability and durability and monetary cost including storage, read, write, and potential migration costs.</li> <li>• Designing light-weight algorithms to optimize data management cost consisting of storage, read, write, potential migration cost across different storage classes based on the status of objects, i.e., hot- and cold-spot.</li> <li>• Designing online and off-line algorithms to make decisions on the number and type of required database instances as well as when on-demand or reserved database instances are deployed.</li> </ul>

*Data consistency and transaction management.* To provide consistent data and support ACID transactions, a concurrency control mechanism should be deployed. This significantly affects the overall performance of response time, availability, and even monetary cost. Thus, it is worth deploying strategies to reduce or exempt coordination across replicas especially across DCs. For this purpose, the concurrency control mechanisms are performed in different levels. The concurrency control in *I/O* level is a widely studied research topic in OLTP applications. The proposed protocols in this level are either lock-free or lock-based, which make a trade-off between response time and memory consumption. **To reduce the need for concurrency in I/O level, it would be relevant (i) to replicate associated data in the servers/DCs that**

are at a close distance, and (ii) order requests against replicas in the network level instead of replication protocol, which has received considerable attention recently. The concurrency mechanisms in the level of *object* and *flow* are appealing since they exempt/reduce the need of coordination. They are in their infancy and currently support a limited number of data objects (e.g., *counters*, some specific *sets* types) and operation types (e.g., increment and decrements for counters, union and intersect for sets). These objects also lack rich relations and operations among themselves that affect reduction in the amount of information maintained and propagated by each replica. The concurrency mechanisms in application and language levels can reduce the need of coordination but they are ad hoc and error-prone.

**Data management cost.** Monetary cost is a key factor for application providers to move their data into storage infrastructure. The optimization of monetary cost is a vital criterion for application providers. The contributing factors in this optimization are: the pricing models, the duration of used resources, the characteristics of workload, and the required QoS depending on the types of applications. We discussed this optimization based on either single/multi QoS metrics. The existing studies dealt with the cost optimization problems in which some cost elements and QoS metrics are considered. Thus, this is a venue to study the cost optimization problems which are complimentary to those already studied. Further work is also needed to investigate when these problems lead to cost trade-offs like storage vs. network, storage vs. cache, and storage vs. computation. Furthermore, due to different classes of storage (for example, S3 Standard, S3 Standard - Infrequent Access, Amazon Glacier, Amazon Reduced Redundancy Storage (RRS)) differentiating in price and performance, it would be interesting to investigate how to replicate objects during their lifetime based on the writes, reads, and size of objects so that the cost is optimized and the required SLA is satisfied. In fact, the objects should be migrated across these classes within or across data stores for this purpose. This mandates to determine the migration time(s) for objects as their status change from hot-spot to cold-spot.

## REFERENCES

- Daniel Abadi. 2012. Consistency Tradeoffs in Modern Distributed Database System Design: CAP is Only Part of the Story. *IEEE Computer* (2012).
- Daniel J. Abadi. 2009. Data Management in the Cloud: Limitations and Opportunities. *Data Engineering Bulletin Issues* (2009).
- Hussam Abu-Libdeh, Lonnie Princehouse, and Hakim Weatherspoon. 2010. RACS: a case for cloud storage diversity (*SoCC '10*).
- Atul Adya. 1999. *Weak Consistency: A Generalized Theory and Optimistic Implementations for Distributed Transactions*. Ph.D. Dissertation. Cambridge, MA, USA.
- Divyakant Agrawal, Amr El Abbadi, Hatem A. Mahmoud, Faisal Nawab, and Kenneth Salem. 2013. *Managing Geo-replicated Data in Multi-datacenters*. Springer Berlin Heidelberg.
- Rakesh Agrawal, Michael J. Carey, and Miron Livny. 1987. Concurrency Control Performance Modeling: Alternatives and Implications. *ACM Transactions on Database Systems* (1987).
- Marcos K. Aguilera, Arif Merchant, Mehul Shah, Alistair Veitch, and Christos Karamanolis. 2007. Sinfonia: A New Paradigm for Building Scalable Distributed Systems. (2007).
- Mustaque Ahamad, Gil Neiger, James E. Burns, Prince Kohli, and Phillip W. Hutto. 1995. Causal memory: definitions, implementation, and programming. *Distributed Computing* (1995).
- Mohammad Al-Fares, Sivasankar Radhakrishnan, Barath Raghavan, Nelson Huang, and Amin Vahdat. 2010. Hedera: Dynamic Flow Scheduling for Data Center Networks (*NSDI '10*).
- Paulo Sérgio Almeida, Ali Shoker, and Carlos Baquero. 2015. Efficient State-Based CRDTs by Delta-Mutation. In *Third International Conference Networked Systems (NETYS)*, 2015.
- Sérgio Almeida, João Leitão, and Luís Rodrigues. 2013. ChainReaction: A Causal+ Consistent Datastore Based on Chain Replication (*EuroSys '13*).
- Peter Alvaro, Peter Bailis, Neil Conway, and Joseph M. Hellerstein. 2013. Consistency Without Borders. In *Proceedings of the 4th Annual Symposium on Cloud Computing (SOCC '13)*.
- Peter Alvaro, Neil Conway, Joseph M. Hellerstein, and David Maier. 2014. Blazes: Coordination analysis for distributed programs. In *International Conference on Data Engineering*, 2014.
- Peter Alvaro, Neil Conway, Joseph M. Hellerstein, and William R. Marczak. 2011. Consistency Analysis in Bloom: a CALM and Collected Approach (*CIDR '11*).
- Eric Anderson, Xiaozhou Li, Mehul A. Shah, Joseph Tucek, and Jay J. Wylie. 2010. What consistency does your key-value store actually provide? (*HotDep '10*).
- Masoud Saeida Ardekani, Pierre Sutra, and Marc Shapiro. 2013. Non-monotonic Snapshot Isolation: Scalable and Strong Consistency for Geo-replicated Transactional Systems (*SRDS '13*).
- Masoud Saeida Ardekani, Pierre Sutra, and Marc Shapiro. 2014. G-DUR: A Middleware for Assembling, Analyzing, and Improving Transactional Protocols (*Middleware '14*).
- Masoud Saeida Ardekani, Pierre Sutra, Marc Shapiro, and Nuno M. Prego. 2013. On the Scalability of Snapshot Isolation (*EuroPar '13*).
- José Enrique Armendáriz-Inigo, A Mauch-Goya, JR de Mendivil, and FD Muñoz-Escófi. 2008. SIPRe: a partial database replication protocol with SI replicas. In *ACM Symposium on Applied Computing (SAC)*, 2008.
- Peter Bailis, Aaron Davidson, Alan Fekete, and et al. 2013. Highly Available Transactions: Virtues and Limitations. *Proceedings of the VLDB Endowment* (2013).

- Peter Bailis, Alan Fekete, Joseph M. Hellerstein, Ali Ghodsi, and Ion Stoica. 2014. Scalable Atomic Visibility with RAMP Transactions (*SIGMOD '14*).
- Peter Bailis, Ali Ghodsi, Joseph M. Hellerstein, and Ion Stoica. 2013. Bolt-on Causal Consistency (*SIGMOD '13*).
- Peter Bailis, Shivaram Venkataraman, Michael J. Franklin, Joseph M. Hellerstein, and Ion Stoica. 2012. Probabilistically bounded staleness for practical partial quorums. *Proceedings of the VLDB Endowment* (2012).
- Jason Baker, Chris Bond, James C. Corbett, and et al. 2011. Megastore: Providing Scalable, Highly Available Storage for Interactive Services (*CIDR '11*).
- Valter Balegas, Sérgio Duarte, Carla Ferreira, and et al. 2015. Putting Consistency Back into Eventual Consistency (*EuroSys '15*).
- Hitesh Ballani, Paolo Costa, Thomas Karagiannis, and Ant Rowstron. 2011. Towards Predictable Datacenter Networks (*SIGCOMM '11*).
- Hal Berenson, Phil Bernstein, Jim Gray, and et al. 1995. A Critique of ANSI SQL Isolation Levels (*SIGMOD '95*).
- David Bermbach, Markus Klems, Stefan Tai, and Michael Menzel. 2011. MetaStorage: A Federated Cloud Storage System to Manage Consistency-Latency Tradeoffs. In *IEEE Conference on Cloud Computing (CLOUD)*, 2011.
- Alysson Bessani, Miguel Correia, Bruno Quaresma, Fernando André, and Paulo Sousa. 2011. DepSky: Dependable and Secure Storage in a Cloud-of-clouds (*EuroSys '11*).
- Carlos Eduardo Bezerra, Fernando Pedone, and Robbert van Renesse. 2014. Scalable State-Machine Replication. In *IEEE/IFIP Conference on Dependable Systems and Networks*, 2014.
- Dhruba Borthakur, Jonathan Gray, Joydeep Sen Sarma, and et al. 2011. Apache Hadoop Goes Realtime at Facebook (*SIGMOD '11*). ACM, New York, NY, USA, 1071–1080.
- Kevin D. Bowers, Ari Juels, and Alina Oprea. 2009. HAIL: A High-availability and Integrity Layer for Cloud Storage. In *ACM Conference on Computer and Communications Security (CCS)*, 2009.
- Matthias Brantner, Daniela Florescu, David Graf, Donald Kossmann, and Tim Kraska. 2008. Building a Database on S3 (*SIGMOD '08*).
- Christian Cachin, Robert Haas, and Marko Vukolic. 2010. Dependable Storage in the Intercloud. In *Research Report RZ*, 2010.
- Brad Calder, Ju Wang, Aaron Ogus, and et al. 2011. Windows Azure Storage: a highly available cloud storage service with strong consistency (*SOSP '11*).
- David G. Campbell, Gopal Kakivaya, and Nigel Ellis. 2010. Extreme Scale with Full SQL Language Support in Microsoft SQL Azure (*SIGMOD '10*).
- Rick Cattell. 2011. Scalable SQL and NoSQL Data Stores. *ACM SIGMOD Record* (2011).
- A. Chan and R. Gray. 1985. Implementing Distributed Read-Only Transactions. *IEEE Transactions on Software Engineering* (1985).
- Chia-Wei Chang, Pangfeng Liu, and Jan-Jan Wu. 2012. Probability-Based Cloud Storage Providers Selection Algorithms with Maximum Availability. In *Proceedings International Conference on Parallel Processing (ICPP)*, 2012.
- Fay Chang, Jeffrey Dean, Sanjay Ghemawat, and et al. 2008. Bigtable: A Distributed Storage System for Structured Data. *ACM Transactions on Computer Systems* (2008).
- Fangfei Chen, Katherine Guo, John Lin, and Thomas La Porta. 2012. Intra-cloud lightning: Building CDNs in the cloud. In *INFOCOM*, 2012.
- H. Chen, H. Jin, and S. Wu. 2016. Minimizing Inter-Server Communications by Exploiting Self-Similarity in Online Social Networks. *IEEE Transactions on Parallel and Distributed Systems* (2016).
- Haopeng Chen, Zhenhua Wang, and Yunmeng Ban. 2013. Access-Load-Aware Dynamic Data Balancing for Cloud Storage Service. In *Conference Internet and Distributed Computing Systems (IDCS)*, 2013.
- Henry C.H. Chen, Yuchong Hu, Patrick P.C. Lee, and Yang Tang. 2014. NCCloud: A Network-Coding-Based Storage System in a Cloud-of-Clouds. *IEEE Trans. Comput.* (2014).
- Kai Chen, Ankit Singla, Atul Singh, and et al. 2014. OSA: An Optical Switching Architecture for Data Center Networks with Unprecedented Flexibility. *IEEE/ACM Transactions on Networking* (2014).
- Houssem-Eddine Chihoub, Shadi Ibrahim, abriel Antoniu, and Maria Perez. 2013. Consistency in the Cloud: When Money Does Matter!. In *IEEE/ACM Symposium on Cluster, Cloud and Grid Computing (CCGrid)*, 2013.
- David Chiu and Gagan Agrawal. 2010. Evaluating caching and storage options on the Amazon Web Services Cloud. In *IEEE/ACM International Conference on Grid Computing*, 2010.
- Asaf Cidon, Robert Escriva, Sachin Katti, Mendel Rosenblum, and Emin Gun Sirer. 2015. Tiered Replication: A Cost-effective Alternative to Full Cluster Geo-replication (*USENIX ATC '15*).
- Asaf Cidon, Stephen M. Rumble, Ryan Stutsman, and et al. 2013. Copysets: Reducing the Frequency of Data Loss in Cloud Storage (*USENIX ATC '13*).
- Brian F. Cooper, Raghu Ramakrishnan, Utkarsh Srivastava, and et al. 2008. PNUTS: Yahoo!'s Hosted Data Serving Platform. *Proceedings of the VLDB Endowment* (2008).
- James C. Corbett, Jeffrey Dean, Michael Epstein, and et al. 2013. Spanner: Google's Globally Distributed Database. *ACM Transactions on Computer Systems* (2013).
- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. 2009. *Introduction to Algorithms* (3. ed.). MIT Press.
- Paolo Costa, Austin Donnelly, Antony Rowstron, and Greg O'Shea. 2012. Camdoop: Exploiting In-network Aggregation for Big Data Applications (*NSDI '12*).
- Anupam Das, Cristian Lumezanu, Yueping Zhang, and et al. 2013b. Transparent and Flexible Network Management for Big Data Processing in the Cloud. In *Presented USENIX Workshop on Hot Topics in Cloud Computing*, 2013.
- Sudipto Das, Divyakant Agrawal, and Amr El Abbadi. 2013a. ElasTraS: An Elastic, Scalable, and Self-managing Transactional Database for the Cloud. *ACM Transactions on Database Systems* (2013).
- Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, and et al. 2007. Dynamo: Amazon's Highly Available Key-value Store (*SOSP '07*).

- Xavier Défago, André Schiper, and Péter Urbán. 2004. Total Order Broadcast and Multicast Algorithms: Taxonomy and Survey. *Comput. Surveys* (2004).
- Alan Demers, Dan Greene, Carl Hauser, and et al. 1987. Epidemic Algorithms for Replicated Database Maintenance (*PODC '87*).
- Alexandros G. Dimakis, P. Brighten Godfrey, Yunnan Wu, Martin J. Wainwright, and Kannan Ramchandran. 2010. Network Coding for Distributed Storage Systems. *IEEE Transactions on Information Theory* (2010).
- Jiaqing Du, Sameh Elnikety, Amitabha Roy, and Willy Zwaenepoel. 2013. Orbe: Scalable Causal Consistency Using Dependency Matrices and Physical Clocks (*SOCC '13*).
- Jiaqing Du, Calin Iorgulescu, Amitabha Roy, and Willy Zwaenepoel. 2014a. Closing The Performance Gap between Causal Consistency and Eventual Consistency. In *Workshop on Principles and Practice of Eventual Consistency (PaPEC)*, 2014.
- Jiaqing Du, Călin Iorgulescu, Amitabha Roy, and Willy Zwaenepoel. 2014b. Gentlerain: Cheap and scalable causal consistency with physical clocks. In *Proceedings of the ACM Symposium on Cloud Computing*, 2014.
- Robert Escriva, Bernard Wong, and Emin Gün Sirer. 2012. HyperDex: A Distributed, Searchable Key-value Store (*SIGCOMM '12*).
- Robert Escriva, Bernard Wong, and Emin Gün Sirer. 2015. Warp: Multi-key transactions for keyvalue stores. (2015).
- Yuan Feng, Baochun Li, and Bo Li. 2012. Postcard: Minimizing Costs on Inter-Datacenter Traffic with Store-and-Forward (*ICDCSW '12*).
- Iliir Fetai and Heiko Schuldt. 2012. Cost-Based Data Consistency in a Data-as-a-Service Cloud Environment. In *IEEE Conference on Cloud Computing (Cloud)*, 2012.
- Daniel Ford, François Labelle, Florentina I Popovici, and et al. 2010. Availability in globally distributed storage systems (*OSDI '10*).
- Seth Gilbert and Nancy Lynch. 2002. Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services. *SIGACT News* (2002).
- Jim Gray and Leslie Lamport. 2006. Consensus on Transaction Commit. *ACM Transactions on Database Systems* (2006).
- Rachid Guerraoui and André Schiper. 2001. Genuine Atomic Multicast in Asynchronous Distributed Systems. *Theoretical Computer Science* (2001).
- Jian Guo, Fangming Liu, Xiaomeng Huang, and et al. 2014. On efficient bandwidth allocation for traffic variability in datacenters. In *IEEE INFOCOM*, 2014.
- R. C. Hansdah and L. M. Patnaik. 1986. *ICDT '86: International Conference on Database Theory Rome, Italy, September 8–10, 1986 Proceedings*. Springer Berlin Heidelberg, Berlin, Heidelberg, Chapter Update serializability in locking, 171–185.
- Zach Hill and Marty Humphrey. 2010. CSAL: A Cloud Storage Abstraction Layer to Enable Portable Cloud Applications. In *Conference on Cloud Computing Technology and Science (CloudCom)*, 2010.
- Cheng Huang, Huseyin Simitci, Yikang Xu, and et al. 2012. Erasure Coding in Windows Azure Storage (*USENIX ATC '12*).
- Anil K. Jain and Richard C. Dubes. 1988. *Algorithms for Clustering Data*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- Sushant Jain, Alok Kumar, Subhasree Mandal, and et al. 2013. B4: Experience with a Globally-deployed Software Defined Wan (*SIGCOMM '13*).
- Vimalkumar Jeyakumar, Mohammad Alizadeh, David Mazières, and et al. 2013. EyeQ: Practical Network Performance Isolation at the Edge (*NSDI'13*).
- Lei Jiao, Jun Li, Tianyin Xu, Wei Du, and Xiaoming Fu. 2016. Optimizing Cost for Online Social Networks on Geo-Distributed Clouds. *IEEE/ACM Transactions on Networking* (2016).
- Joarder Kamal, Manzur Murshed, and Rajkumar Buyya. 2016. Workload-aware Incremental Repartitioning of Shared-nothing Distributed Databases for Scalable OLTP Applications. *Future Generation Computer Systems* (2016).
- Osama Khan, Randal C Burns, James S Plank, William Pierce, and Cheng Huang. 2012. Rethinking erasure codes for cloud file systems: minimizing I/O for recovery and degraded reads. (*FAST '12*).
- Tadeusz Kobus, Maciej Kokocinski, and Pawel T. Wojciechowski. 2013. Hybrid Replication: State-Machine-Based and Deferred-Update Replication Schemes Combined. In *ICDCS*, 2013.
- Ramakrishna Kotla, Lorenzo Alvisi, and Mike Dahlin. 2007. SafeStore: A Durable and Practical Storage System. In *2007 USENIX Annual Technical Conference on Proceedings of the USENIX Annual Technical Conference (ATC'07)*. USENIX Association, Berkeley, CA, USA, Article 10, 14 pages.
- Tim Kraska, Martin Hentschel, Gustavo Alonso, and Donald Kossmann. 2009. Consistency Rationing in the Cloud: Pay Only when It Matters. *Proceedings of the VLDB Endowment* (2009).
- Diego Kreutz, Fernando M. V. Ramos, Paulo Verissimo, and et al. 2015. Software-Defined Networking: A Comprehensive Survey. *Proc. IEEE* (2015).
- K. Ashwin Kumar, Abdul Quamar, Amol Deshpande, and Samir Khuller. 2014. SWORD: workload-aware data placement and replica selection for cloud data management systems. *The VLDB Journal* (2014).
- H. T. Kung and John T. Robinson. 1981. On Optimistic Methods for Concurrency Control. *ACM Transactions on Database Systems* (1981).
- Avinash Lakshman and Prashant Malik. 2010. Cassandra: a decentralized structured storage system. *ACM SIGOPS Operating Systems Review* (2010).
- Leslie Lamport. 1998. The Part-time Parliament. *ACM Transactions on Computer Systems* (1998).
- Leslie Lamport. 2005. *Generalized Consensus and Paxos*. Technical Report MSR-TR-2005-33.
- Leslie Lamport. 2006. Fast Paxos. *Distributed Computing* 19, 2 (2006), 79–103.
- Jeongkeun Lee, Yoshio Turner, Myungjin Lee, and et al. 2014. Application-driven Bandwidth Guarantees in Datacenters (*SIGCOMM '14*).
- Cheng Li, João Leitão, Allen Clement, and et al. 2014. Automating the Choice of Consistency Levels in Replicated Systems (*USENIX ATC '14*).
- Cheng Li, Daniel Porto, Allen Clement, and et al. 2012. Making Geo-replicated Systems Fast As Possible, Consistent when Necessary (*OSDI '12*).



- Jialin Li, Ellis Michael, Naveen Kr. Sharma, Adriana Szekeres, and Dan R. K. Ports. 2016. Just Say NO to Paxos Overhead: Replacing Consensus with Network Ordering. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. USENIX Association, GA, 467–483.
- Mingqiang Li, Chuan Qin, and Patrick P. C. Lee. 2015. CDStore: Toward Reliable, Secure, and Cost-Efficient Cloud Storage via Convergent Dispersal (*USENIX ATC '15*).
- Wenhao Li, Yun Yang, Jinjun Chen, and Dong Yuan. 2012. A Cost-Effective Mechanism for Cloud Data Reliability Management Based on Proactive Replica Checking. In *Proceedings IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID), 2012*.
- Wenhao Li, Yun Yang, and Dong Yuan. 2011. A Novel Cost-Effective Dynamic Data Replication Strategy for Reliability in Cloud Data Centres. In *IEEE Conference on Dependable, Autonomic and Secure Computing (DASC), 2011*.
- Guanfeng Liang and Ulaş C. Kozat. 2014. Fast Cloud: Pushing the Envelope on Delay Performance of Cloud Storage with Coding. *IEEE/ACM Transactions on Networking* (2014).
- Guanfeng Liang and Ula C. Kozat. 2015. On Throughput-Delay Optimal Access to Storage Clouds via Load Adaptive Coding and Chunking. *IEEE/ACM Transactions on Networking* (2015).
- Guoxin Liu, Haiying Shen, and Harrison Chandler. 2013. Selective Data replication for Online Social Networks with Distributed Datacenters. In *IEEE International Conference on Network Protocols (ICNP), 2013*.
- J. Liu and H. Shen. 2016. A Low-Cost Multi-failure Resilient Replication Scheme for High Data Availability in Cloud Storage. In *IEEE 23rd International Conference on High Performance Computing (HiPC)*. 242–251. DOI: <http://dx.doi.org/10.1109/HiPC.2016.036>
- Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Andersen. 2011. Don'T Settle for Eventual: Scalable Causal Consistency for Wide-area Storage with COPS (*SOSP '11*).
- Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Andersen. 2013. Stronger Semantics for Low-latency Geo-replicated Storage (*NSDI '13*).
- David Lomet. 1996. Replicated Indexes for Distributed Data. In *Conference on Parallel and Distributed Information Systems, 1996*.
- Yadi Ma, Thyaga Nandagopal, Krishna PN Puttaswamy, and Suman Banerjee. 2013. An Ensemble of Replication and Erasure Codes for Cloud File Systems. In *INFOCOM, 2013*.
- Hatem Mahmoud, Faisal Nawab, Alexander Pucher, Divyakant Agrawal, and Amr El Abbadi. 2013. Low-latency Multi-datacenter Databases Using Replicated Commit. *Proceedings of the VLDB Endowment* (2013).
- Hatem A. Mahmoud, Vaibhav Arora, Faisal Nawab, Divyakant Agrawal, and Amr El Abbadi. 2014. MaaT: Effective and Scalable Coordination of Distributed Transactions in the Cloud. *Proceedings of the VLDB Endowment* (2014).
- Yaser Mansouri and Rajkumar Buyya. 2016. To Move or Not to Move: Cost Optimization in a Dual Cloud-based Storage Architecture. *Journal of Network and Computer Applications* (2016).
- Yaser Mansouri, Adel Nadjaran Toosi, and Rajkumar Buyya. 2013. Brokering Algorithms for Optimizing the Availability and Cost of Cloud Storage Services. In *Conference on Cloud Computing Technology and Science (Cloudcom), 2013*.
- Y. Mansouri, A. Nadjaran Toosi, and R. Buyya. 2017. Cost Optimization for Dynamic Replication and Migration of Data in Cloud Data Centers. *IEEE Transactions on Cloud Computing* (2017). DOI: <http://dx.doi.org/10.1109/TCC.2017.2659728>
- Bo Mao, Suzhen Wu, and Hong Jiang. 2016. Exploiting Workload Characteristics and Service Diversity to Improve the Availability of Cloud Storage Systems. *IEEE Transaction Parallel Distributed Systems* 27, 7 (2016), 2010–2021.
- John C. McCullough, John Dunagan, Alec Wolman, and Alex C. Snoeren. 2010. Stout: An Adaptive Interface to Scalable Cloud Storage (*USENIXATC '10*).
- Iulian Moraru, David G. Andersen, and Michael Kaminsky. 2013. There is More Consensus in Egalitarian Parliaments. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles (SOSP '13)*. ACM, New York, NY, USA, 358–372.
- Shuai Mu, Yang Cui, Yang Zhang, Wyatt Lloyd, and Jinyang Li. 2014. Extracting More Concurrency from Distributed Transactions (*OSDI '14*).
- Subramanian Muralidhar, Wyatt Lloyd, Sabyasachi Roy, and et al. 2014. f4: Facebook's Warm BLOB Storage System (*OSDI '14*).
- Maurizio Naldi and Loretta Mastroeni. 2013. Cloud Storage Pricing: A Comparison of Current Practices. In *Proceedings of the 2013 International Workshop on Hot Topics in Cloud Services (HotTopICS '13)*.
- Rajesh Nishtala, Hans Fugal, Steven Grimm, and et al. 2013. Scaling Memcache at Facebook (*NSDI '13*).
- Brian M. Oki and Barbara H. Liskov. 1988. Viewstamped Replication: A New Primary Copy Method to Support Highly-Available Distributed Systems (*PODC '88*).
- Leandro Pacheco, Daniele Sciascia, and Fernando Pedone. 2014. Parallel Deferred Update Replication. In *IEEE Symposium on Network Computing and Applications (NCA), 2014*.
- Fernando Pedone, Matthias Wiesmann, André Schiper, Bettina Kemme, and Gustavo Alonso. 2000. Understanding Replication in Databases and Distributed Systems. In *ICDCS, 2000*.
- Sebastiano Peluso, Pedro Ruivo, Paolo Romano, Francesco Quaglia, and Luis Rodrigues. 2012. When Scalability Meets Consistency: Genuine Multiversion Update-Serializable Partial Data Replication. In *ICDCS, 2012*.
- Lucian Popa, Praveen Yalagandula, Sujata Banerjee, and et al. 2013. ElasticSwitch: Practical Work-conserving Bandwidth Guarantees for Cloud Computing (*SIGCOMM '13*).
- Dan R. K. Ports, Jialin Li, Vincent Liu, Naveen Kr. Sharma, and Arvind Krishnamurthy. 2015. Designing Distributed Systems Using Approximate Synchrony in Data Center Networks. In *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI 15)*. USENIX Association, Oakland, CA, 43–57.
- Josep M. Pujol, Vijay Erramilli, Georgos Siganos, and et al. 2010. The Little Engine(s) That Could: Scaling Online Social Networks (*SIGCOMM '10*).
- Krishna P.N. Puttaswamy, Thyaga Nandagopal, and Murali Kodialam. 2012. Frugal Storage for Cloud File Systems (*EuroSys '12*).
- K.V. Rashmi, N.B. Shah, and P.V. Kumar. 2011. Optimal Exact-Regenerating Codes for Distributed Storage at the MSR and MBR Points via a Product-Matrix Construction. *IEEE Transactions on Information Theory* (2011).



- K.V. Rashmi, Nihar B. Shah, Dikang Gu, and et al. 2014. A "Hitchhiker's" Guide to Fast and Efficient Data Reconstruction in Erasure-coded Data Centers (*SIGCOMM '14*).
- Ron Roth. 2006. *Introduction to Coding Theory*. Cambridge University Press, 2006.
- Sudip Roy, Lucja Kot, Gabriel Bender, and et al. 2015. The Homeostasis Protocol: Avoiding Transaction Coordination Through Program Analysis (*SIGMOD '15*).
- Stephen M. Rumble, Ankita Kejriwal, and John Ousterhout. 2014. Log-structured Memory for DRAM-based Storage. In *Proceedings of the 12th USENIX Conference on File and Storage Technologies (FAST'14)*.
- Yasushi Saito and Marc Shapiro. 2005. Optimistic Replication. *ACM Comput. Surveys* (2005).
- Sherif Sakr. 2014. Cloud-hosted databases: technologies, challenges and opportunities. *Cluster Computing* 17, 2 (01 Jun 2014), 487–502.
- S. Sakr, A. Liu, D.M. Batista, and M. Alomari. 2011. A Survey of Large Scale Data Management Approaches in Cloud Environments. *IEEE Communications Surveys Tutorials*, (2011).
- Maheswaran Sathiamoorthy, Megasthenis Asteris, Dimitris Papailiopoulos, and et al. 2013. XORing Elephants: Novel Erasure Codes for Big Data. *Proceedings of the VLDB Endowment* (2013).
- Nicolas Schiper, Pierre Sutra, and Fernando Pedone. 2009. Genuine Versus Non-Genuine Atomic Multicast Protocols for Wide Area Networks: An Empirical Study (*SRDS '09*).
- Nicolas Schiper, Pierre Sutra, and Fernando Pedone. 2010. P-Store: Genuine Partial Replication in Wide Area Networks (*SRDS '10*).
- Daniele Sciascia, Fernando Pedone, and Flavio Junqueira. 2012. Scalable Deferred Update Replication (*DSN '12*).
- Steven S. Seiden. 2000. A Guessing Game and Randomized Online Algorithms (*STOC '00*).
- Nihar B. Shah, Kangwook Lee, and Kannan Ramchandran. 2014. The MDS queue: Analysing the latency performance of erasure codes. In *International Symposium on Information Theory, 2014*.
- Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. 2011. Conflict-free Replicated Data Types. In *International Conference on Stabilization, Safety, and Security of Distributed Systems, 2011*.
- Artyom Sharov, Alexander Shraer, Arif Merchant, and Murray Stokely. 2015. Take Me to Your Leader!: Online Optimization of Distributed Storage Configurations. *Proceedings of the VLDB Endowment* (2015).
- Dennis Shasha, Francois Llirbat, Eric Simon, and Patrick Valduriez. 1995. Transaction Chopping: Algorithms and Performance Studies. *ACM Transactions on Database Systems* (1995).
- Min Shen, Ajay D. Kshemkalyani, and Ta Yuan Hsu. 2015. OPCAM: Optimal Algorithms Implementing Causal Memories in Shared Memory Systems. In *International Conference on Distributed Computing and Networking (ICDCN), 2015*.
- Alan Shieh, Srikanth Kandula, Albert Greenberg, Changhoon Kim, and Bikas Saha. 2011. Sharing the Data Center Network (*NSDI '11*).
- Youngjoo Shin, Dongyoung Koo, and Junbeom Hur. 2017. A Survey of Secure Data Deduplication Schemes for Cloud Storage Systems. *ACM Comput. Surv.* 49, 4, Article 74 (Jan. 2017), 38 pages.
- David Shue, Michael J. Freedman, and Anees Shaikh. 2012. Performance Isolation and Fairness for Multi-tenant Cloud Storage (*OSDI '12*).
- Jeff Shute, Radek Vingralek, Bart Samwel, and et al. 2013. F1: A Distributed SQL Database That Scales. *Proceedings of the VLDB Endowment* (2013).
- KC Sivaramakrishnan, Gowtham Kaki, and Suresh Jagannathan. 2015. Declarative Programming over Eventually Consistent Data Stores (*PLDI '15*).
- Yair Sovran, Russell Power, Marcos K. Aguilera, and Jinyang Li. 2011. Transactional Storage for Geo-replicated Systems (*SOSP '11*).
- Josef Spillner, Gerd Bombach, Steffen Matthischke, and et al. 2011. Information Dispersion over Redundant Arrays of Optimal Cloud Storage for Desktop Users. In *IEEE Conference on Utility and Cloud Computing (UCC), 2011*.
- Christopher Stewart, Aniket Chakrabarti, and Rean Griffith. 2013. Zoolander: Efficiently Meeting Very Strict, Low-Latency SLOs. In *International Conference on Autonomic Computing (ICAC), 2013*.
- Changho Suh and Kannan Ramchandran. 2011. Exact-Repair MDS Code Construction Using Interference Alignment. *IEEE Transactions on Information Theory* (2011).
- Lalith Suresh, Marco Canini, Stefan Schmid, and Anja Feldmann. 2015. C3: Cutting Tail Latency in Cloud Data Stores via Adaptive Replica Selection (*NSDI '15*).
- Andrew S. Tanenbaum and Maarten van Steen. 2006. *Distributed Systems: Principles and Paradigms (2Nd Edition)*. Prentice-Hall, Inc.
- Jing Tang, Xueyan Tang, and Junsong Yuan. 2015. Optimizing inter-server communication for online social networks. In *ICDCS, 2015*.
- Jeff Terrace and Michael J. Freedman. 2009. Object Storage on CRAQ: High-throughput Chain Replication for Read-mostly Workloads (*USENIX '09*).
- Douglas B. Terry, Vijayan Prabhakaran, Ramakrishna Kotla, and et al. 2013. Consistency-based service level agreements for cloud storage (*SOSP '13*).
- Robert H. Thomas. 1979. A Majority Consensus Approach to Concurrency Control for Multiple Copy Databases. *ACM Transactions on Database Systems* (1979).
- Alexander Thomson, Thaddeus Diamond, Shu-Chun Weng, and et al. 2012. Calvin: Fast Distributed Transactions for Partitioned Database Systems (*SIGMOD '12*).
- Ashish Thusoo, Joydeep Sen Sarma, Namit Jain, and et al. 2010. Hive - a petabyte scale data warehouse using Hadoop. In *International Conference on Data Engineering (ICDE), 2010*.
- Duc A. Tran, Khanh Nguyen, and Cuong Pham. 2012. S-CLONE: Socially-aware data replication for social networks. *Computer Networks* (2012).

- Nguyen Tran, Marcos K. Aguilera, and Mahesh Balakrishnan. 2011. Online Migration for Geo-distributed Storage Systems (*USENIXATC '11*).
- Balajee Vamanan, Jahangir Hasan, and T.N. Vijaykumar. 2012. Deadline-aware Datacenter TCP (D2TCP) (*SIGCOMM '12*).
- Robbert van Renesse and Fred B. Schneider. 2004. Chain Replication for Supporting High Throughput and Availability (*OSDI '04*).
- Srikumar Venugopal, Rajkumar Buyya, and Kotagiri Ramamohanarao. 2006. A Taxonomy of Data Grids for Distributed Data Sharing, Management, and Processing. *ACM Comput. Surv.* 38, 1, Article 3 (June 2006).
- Hoang Tam Vo, Sheng Wang, Divyakant Agrawal, Gang Chen, and Beng Chin Ooi. 2012. LogBase: A Scalable Log-structured Database System in the Cloud. *Proceedings of the VLDB Endowment* (2012).
- Hiroshi Wada, Alan Fekete, Liang Zhao, Kevin Lee, and Anna Liu. 2011. Data Consistency Properties and the Trade-offs in Commercial Cloud Storage: the Consumers' Perspective. In *Conference on Innovative Data Systems (CIDR)*, 2011.
- Ting Wang, Zhiyang Su, Yu Xia, and Mounir Hamdi. 2014. Rethinking the Data Center Networking: Architecture, Network Protocols, and Resource Sharing. *IEEE Access* (2014).
- Wei Wang, Baochun Li, and Ben Liang. 2013. To Reserve or Not to Reserve: Optimal Online Multi-Instance Acquisition in IaaS Clouds. In *International Conference on Autonomic Computing (ICAC)*, 2013.
- Yunnan Wu, Alexandros G Dimakis, and Kannan Ramchandran. 2007. Deterministic regenerating codes for distributed storage. *Allerton Conference on Control, Computing, and Communication* (2007).
- Yu Wu, Zhizhong Zhang, Chuan Wu, and et al. 2015b. Orchestrating Bulk Data Transfers across Geo-Distributed Datacenters. *IEEE Transactions on Cloud Computing* (2015).
- Zhe Wu, Michael Butkiewicz, Dorian Perkins, Ethan Katz-Bassett, and Harsha V. Madhyastha. 2013. SPANStore: Cost-effective Geo-replicated Storage Spanning Multiple Cloud Services (*SOSP '13*).
- Zhe Wu, Curtis Yu, and Harsha V. Madhyastha. 2015a. CosTLO: Cost-effective Redundancy for Lower Latency Variance on Cloud Storage Services (*NSDI '15*).
- Mingyuan Xia, Mohit Saxena, Mario Blaum, and David A. Pease. 2015. A Tale of Two Erasure Codes in HDFS (*FAST '15*).
- Yu Xiang, Tian Lan, Vaneet Aggarwal, and Yih Farn R. Chen. 2014. Joint Latency and Cost Optimization for Erasure-coded Data Center Storage. *ACM SIGMETRICS Performance Evaluation Review* (2014).
- Di Xie, Ning Ding, Y. Charlie Hu, and Ramana Kompella. 2012. The Only Constant is Change: Incorporating Time-varying Network Reservations in Data Centers (*SIGCOMM '12*).
- Boyang Yu and Jianping Pan. 2015. Location-aware associated data placement for geo-distributed data-intensive applications. In *INFOCOM*, 2015.
- Haifeng Yu and Amin Vahdat. 2002. Design and evaluation of a conit-based continuous consistency model for replicated services. *ACM Transaction Computer Systems* (2002).
- Wenying Zeng, Yuelong Zhao, Kairi Ou, and Wei Song. 2009. Research on cloud storage architecture and key technologies (*ICIS '09*).
- Hong Zhang, Kai Chen, Wei Bai, and et al. 2015. Guaranteeing Deadlines for Inter-datacenter Transfers (*EuroSys '15*).
- Irene Zhang, Naveen Kr. Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan R. K. Ports. 2015. Building Consistent Transactions with Inconsistent Replication (*SOSP '15*).
- Quanlu Zhang, Shenglong Li, Zhenhua Li, and et al. 2015. CHARM: A Cost-efficient Multi-cloud Data Hosting Scheme with High Availability. *IEEE Transactions on Cloud Computing* (2015).
- Xuyun Zhang, Chang Liu, Surya Nepal, Suraj Pandey, and Jinjun Chen. 2013a. A Privacy Leakage Upper Bound Constraint-Based Approach for Cost-Effective Privacy Preserving of Intermediate Data Sets in Cloud. *IEEE Transactions on Parallel and Distributed Systems* (2013).
- Yang Zhang, Russell Power, Siyuan Zhou, and et al. 2013b. Transaction Chains: Achieving Serializability with Low Latency in Geo-distributed Storage Systems (*SOSP '13*).
- Zhe Zhang, Amey Deshpande, Xiaosong Ma, Eno Thereska, and Dushyanth Narayanan. 2010. Does erasure coding have a role to play in my data center. *Microsoft research MSR-TR-2010* (2010).
- L. Zhao, S. Sakr, and A. Liu. 2015. A Framework for Consumer-Centric SLA Management of Cloud-Hosted Databases. *IEEE Transactions on Services Computing* (July 2015).
- Jingya Zhou, Jianxi Fan, Jin Wang, Baolei Cheng, and Juncheng Jia. 2016. Towards traffic minimization for data placement in online social networks. *Concurrency and Computation: Practice and Experience* (2016).

## Online Appendix to: Data Storage Management in Cloud Environments: Taxonomy, Survey, and Future Directions

Yaser Mansouri, Adel Nadjaran Toosi,  
and Rajkumar Buyya, The University of Melbourne, Australia.

This Appendix consists of 7 sub Appendixes/tables. For concise, we use the following abbreviations in the appendix: DC (Data Center), SDC (Single Data Center), MDC (Multi-Data Center), KV (Key Value), RDB (Relational DataBase), RO (Read-only), WO (Write-Only), RW (Read-Wrire), CL (Commit Latency), and n/a (Not Applicable).

The projects listed in Tables I and II are either open-source or commercial. The reference of these projects is as follows:

BigTable: <https://cloud.google.com/bigtable/>

S3: <https://aws.amazon.com/s3/>

SimpleDB: <https://aws.amazon.com/simpliedb/>

AmazonRDS: <https://aws.amazon.com/rds/>

SQL Azure: <https://azure.microsoft.com/en-us/services/sql-database/>

Azure Blob: <https://azure.microsoft.com/en-au/services/storage/blobs/>

Cassandra: <http://cassandra.apache.org/>

MongoDB: <https://www.mongodb.com/>

Volemort: <http://www.project-voldemort.com/voldemort/>

Riak: <http://basho.com/products/#riak>

Hbase: <https://hbase.apache.org/>

Google Cloud Datastore: <https://cloud.google.com/datastore/>

Google Cloud SQL: <https://cloud.google.com/sql/>

VoltDB: <https://www.voltdb.com/>

Window Azure Storage (LRS, ZRS, GRS, RA-GRS): <https://docs.microsoft.com/en-us/azure/storage/common/storage-redundancy>

## A. CHARACTERISTICS OF DIFFERENT DATA STORE

Table I: Summary of different data stores

Data store	Good for what?	Data Structure	Data Model	Replication mechanism	Consistency level	CAP option	PACELC option	Transaction support	License
BigTable*	OLAP	(Un-/)structured	Extensible record	Single master (lazy)	Eventual	CP	PC/EC	NO	Internal@Google
PNUTS	OLAP	Structured	key-value	Single master (lazy)	Ordering/Eventual	CP	PC/EL	NO	Internal@Yahoo
Dynamo**	OLAP	Un-structured	Key-value	Quorum	Eventual	AP	PA/EL	NO	Internal@Amazon
S3***	OLAP	Un-structured	Key-value	Multi-master (lazy)	Eventual	AP	PA/EL	NO	Commercial
SimpleDB†	OLAP	(Semi-/)structured	Document	Single master (lazy)	Eventual†	AP	PA/EC	NO	Commercial
Amazon RDS	OLTP	Structured	Relational	Single master (eager/lazy)	Strong	CA	PC/EC	YES	Commercial
SQL Azure††	OLTP	Structured	Relational	Multi-master (Lazy)	Strong	CA	PC/EC	YES	Commercial
Azure Blob† † †	OLAP	Un-structured	Key-value	n/a (Eager)	Eventual/Strong	CP	PC/EC	NO	Commercial
Cassandra	OLAP	Un-structured	Extensible record	Quorum	Tunable◊	AP	PA/EL	NO	Open source
CouchDB	OLAP	Semi-structured	Document	Multi-master (lazy)	Eventual	AP	PA/EL	NO	Open source
MongoDB	OLAP	Semi-structured	Document	Single master (lazy)	Eventual	AP	PA/EC	NO	Open source
Voldemort	OLAP	Un-structured	Key-value	Quorum	Eventual/Strict quorum	AP	PA/EL	NO	Internal@LinkedIn
Riak	OLAP	Un-structured	key-value	Quorum	Tunable	AP	PA/EL	NO	Open source
HBase	OLAP	(Semi-/)structured	Extensible record	Single master (lazy)	Strong◊◊	CP	PC/EC	NO	Open source
Google Cloud Datastore‡	OLAP	(Semi/Un)-structured	Extensible record	Multi-master (eager)	Eventual/Strong	CP	PC/EC	YES	Commercial
Google Cloud SQL‡‡	OLTP	Structured	Relational	Single master (n/a)	Strong	AC	PC/EL	YES	Commercial
VoltDB‡ † †	OLTP	Structured	Relational	Single/Multi-master (eager)	Strong	CP	PC/EC	YES	Commercial (VoltDB, Inc.)

\* HyperTable is an open source BigTable implementation. Spanner [Corbett et al. 2013] is layered on a implementation of BigTable to provide fully relational databases, and F1 [Shute et al. 2013] was built based on Spanner to implement a relational database based on MySQL.

\*\* Amazon DynamoDB is built on the principles of Dynamo and supports eventual consistency by default and strong consistency (read-after-write) via tuning the parameters of quorum replication (see Section 4.1.5). \*\*\* Google cloud storage has characteristics same as S3. S3 provides read-after-write consistency for puts and eventual consistency for update and delete operations. This replication mechanism is for cross-region replication.

† SimpleDB supports two consistency semantic: eventual consistency read as defined in Section 5.4.1 and consistent read in which the returning results reflect writes that received a successful response prior to the read.

†† SQL Azure supports multi-master with lazy propagation data replication across regions. It also provides Geo-replication in which data replication mechanism is single-master with lazy propagation.

††† Azure blob storage supports Locally/Zone redundant storage (LRS/ZRS), Geo-redundant storage (GRS), and Read-access geo-redundant storage (RA-GRS) storage. In LRS, data is replicated within a DC with eager propagation, while in ZRS, GRS, and RA-GRS data is replicated across DCs with the help of a single master replica with lazy propagation. Azure Table is the same as Azure blob in all aspects, but it stores structured data.

◊ Cassandra offers *tunable consistency* as an extension of eventual consistency by which an application provider decides a write requests must be initially directed to which replica(s) to update data based on its requirements like availability and latency.

◊◊ Hbase supports strong consistency by default if data is retrieved from master replica; otherwise consistency semantic is eventual. It also supports lightweight transactions (e.g., set and compare) not distributed transactions.

‡ Google Cloud Data provides *strong consistency* by default for *ancestor* queries executing against an entity group (see sections 2.2 and 3.3), and *eventual consistency* for *global* queries executing against entity groups.

‡‡ Google Cloud SQL allows a read replica as a copy of master replica in any zone. Read replica reflects the changes to master replica in almost real time. This data store handles replica failover as master replica outages. Hence, since it gives priority to low latency (i.e., high availability) over strong consistency, we call it as a PC/EL system; otherwise it is a PC/EC system.

‡ † † VoltDB is a commercial version of H-store which was shut down in 2016.

## B. CONSISTENCY VERSUS LATENCY IN DIFFERENT REPLICATION STRATEGIES

Table II: Consistency-latency trade-off of different replication techniques.

Update initiation	Update propagation	Latency source	What is tradeoff?	Consistency semantic	Examples of data stores
Single master	Eager	<b>Case 1:</b> The latency within DCs (very small) and across DCs.	Trading latency for consistency	Strong	Amazon RDS $\diamond$
	Lazy	<b>Case 2:</b> The latency consists of routing reads to master replica and waiting time in queue for serving reads and recovery time for probable failed master replica.	Trading consistency for latency and vice versa	Strong	PNUTS, MongoDB
Multi-master	Eager	<b>Case 3:</b> The latency is the same as in Case 1 plus conflict detection and resolution time.	Trading consistency for latency	Strong	Window Azure Storage (LRS) $\diamond\diamond$
	Lazy	<b>Case 4:</b> The latency is the same as that in case 2.	The same Trade-off as in Case 2	Eventual	S3 (with Geo-replication setting)
Quorum	Eager	<b>Case 5:</b> The latency is similar to that in Case 2.	Trading consistency for latency	Conditional strong $\dagger$	Dynamo, Riak, Cassandra*
	Lazy	<b>Case 6:</b> The latency is the same as that in Case 2.	The same Trade-off as in Case 2	Conditional strong $\dagger\dagger$	

$\dagger$  If reads are submitted to synchronous updated replica precipitated in quorum, then consistency is guaranteed but latency is high.

$\dagger\dagger$  If reads are submitted to asynchronous updated replica precipitated in quorum, then consistency is not guaranteed.

$\diamond$  Amazon RDS provides two distinct replication options: Multi-AZ Deployments in which objects are eagerly updated and Read Replica in which objects are lazily updated.

$\diamond\diamond$  As discussed in Appendix A, Windows Azure Storage (WAS) offers four replication options: LRS, ZRS, GRS, and RA-GRS. Apart from LRS, all options support data replication across DCs, but only RA-GRS allows to read data from a secondary DC as well as the primary DC. Moreover, LRS option provides data replication within a DC with eager propagation, while other options support data replication with lazy propagation across DCs.

\* These data stores use a combination of both update propagation techniques.

### C. PROJECTS WITH WEAK CONSISTENCY LEVEL

Table III: Summary of Projects with Weak Consistency Level across key-value data stores. All these projects use Multi-master replication with lazy propagation.

Project specifications	Replication mechanism	Consistency level	Conflict resolution	Transaction support	RL <sup>†</sup>	Fault tolerance	Features
<b>COPS</b> [Lloyd et al. 2011]	<ul style="list-style-type: none"> <li>Full</li> <li>Multi-row</li> </ul>	Causal+	LWW pre-defined procedure	RO	$\leq 4$	Client Server DC	<ul style="list-style-type: none"> <li>reduces the number of dependency checks by maintaining the nearest dependencies for each value.</li> <li>is not scalable if applications interact with DC for a long session.</li> </ul>
<b>Eiger</b> <sup>††</sup> [Lloyd et al. 2013]	<ul style="list-style-type: none"> <li>Full</li> <li>Multi-row</li> </ul>	Causal+	LWW pre-defined procedure	RO WO	$\leq 6$	Client Server DC	<ul style="list-style-type: none"> <li>is more scalable (due to storing fewer dependencies) and efficient (due to using dependencies on operations rather than values) compared to COPS.</li> </ul>
<b>ChainReaction</b> [Almeida et al. 2013]	<ul style="list-style-type: none"> <li>Full</li> <li>Multi-row</li> </ul>	Causal+	LWW	RO	$\geq 4$	Server	<ul style="list-style-type: none"> <li>maintains and exchanges less metadata than COPS.</li> </ul>
<b>Orbe</b> [Du et al. 2013b]	<ul style="list-style-type: none"> <li>Full</li> <li>Multi-row</li> </ul>	Causal+	Pre-defined procedure	RO	4	DC	<ul style="list-style-type: none"> <li>reduces the size of dependencies metadata by means of causality transitivity and sparse matrix.</li> <li>supports RO transactions by using dependencies matrix and physical clock</li> </ul>
<b>Shim</b> [Bailis et al. 2013]	<ul style="list-style-type: none"> <li>Full</li> <li>Multi-row</li> </ul>	Causal+	LWW	RO WO	$\geq 2$	n/a	<ul style="list-style-type: none"> <li>is a layer between applications and eventually consistent data stores to guarantee causal+.</li> <li>can switch between eventual and causal consistency.</li> </ul>
<b>GentleRain</b> [Du et al. 2014a]	<ul style="list-style-type: none"> <li>Full</li> <li>Multi-row</li> </ul>	Causal+	Pre-defined procedure	RO $\diamond$	2	×	<ul style="list-style-type: none"> <li>achieves higher throughput and scalability compared to COPS and Orbe via ignoring dependency check messages for write operations at the expense of more latency for them.</li> </ul>
<b>Opt-Track and Opt-Track-CRP</b> [Shen et al. 2015]	<ul style="list-style-type: none"> <li>Partial, full</li> <li>Multi-row</li> </ul>	Causal	n/a	n/a	n/a	×	<ul style="list-style-type: none"> <li>uses full replication when read operations are issued from almost all DCs in the system.</li> </ul>
<b>Hyksos</b> [Nawab et al. 2015]	<ul style="list-style-type: none"> <li>Full</li> <li>Multi-row</li> </ul>	Causal	n/a	RO	$\geq 2$	n/a	<ul style="list-style-type: none"> <li>uses a <i>distributed shared log</i> to serialize operations rather than a centralized one as used in Message Future [Nawab et al. 2013], Megastore [Baker et al. 2011], and Paxos-CP [Patterson et al. 2012].</li> </ul>

<sup>†</sup> Read Latency (RL) is measured in the number of messages.

<sup>††</sup> It supports *columin-family* data model.

$\diamond$  Read-snapshot transactions are the same as read transactions but read snapshot transactions must guarantee that the snapshot includes any values previously seen by the client. They are more efficient than RO transactions.

## D. PROJECTS WITH ADAPTIVE CONSISTENCY LEVEL

Table IV: Summary of Projects with Adaptive Consistency Level **within and across DCs**.

Project specifications	Replication mechanism	Consistency levels	Conflict solution	Tra.† type	Trade-off	Features
<b>PRACTI</b> [Belaramani et al. 2006] SDC	<ul style="list-style-type: none"> <li>Partial, multi-shard</li> <li>Multi-master, lazy/eager</li> </ul>	<ul style="list-style-type: none"> <li>Causal</li> <li>Ordering</li> <li>SSER</li> </ul>	Pre-defined procedure	n/a	Availability-consistency	<ul style="list-style-type: none"> <li>is not fault-tolerant.</li> </ul>
<b>Brantner et al. [2008]</b> SDC	<ul style="list-style-type: none"> <li>Full, single row</li> <li>n/a, n/a</li> </ul>	<ul style="list-style-type: none"> <li>Eventual</li> <li>RYW††</li> </ul>	n/a	n/a	Cost-consistency	<ul style="list-style-type: none"> <li>builds protocols on S3.</li> <li>handles client failures.</li> </ul>
<b>CRAQ</b> [Terrace and Freedman 2009] SDC	<ul style="list-style-type: none"> <li>Full, single row</li> <li>Single master, lazy</li> </ul>	<ul style="list-style-type: none"> <li>Eventual</li> <li>Strong</li> </ul>	n/a	Mini-Tran.*	Availability-consistency	<ul style="list-style-type: none"> <li>provides eventual consistency with maximum bounded inconsistency window in terms of version and time.</li> <li>tolerates node failures.</li> </ul>
<b>Consistency rationing</b> [Kraska et al. 2009] SDC	<ul style="list-style-type: none"> <li>Full, single shard</li> <li>n/a, n/a</li> </ul>	<ul style="list-style-type: none"> <li>Session</li> <li>SI</li> </ul>	LWW	RW	Cost-consistency	<ul style="list-style-type: none"> <li>provides consistency levels on top of S3.</li> <li>does not handle any types of failures.</li> </ul>
<b>Bernbach et al. [2011]</b> SDC	<ul style="list-style-type: none"> <li>Full, single row</li> <li>n/a, n/a</li> </ul>	<ul style="list-style-type: none"> <li>MR</li> <li>RYW</li> </ul>	Pre-defined procedure	n/a	n/a	<ul style="list-style-type: none"> <li>inserts a middle-ware between client and underneath S3, DynamoDB, and SimpleDB.</li> <li>is not fault-tolerant.</li> </ul>
<b>C3</b> [Fetai and Schuldt 2012] SDC	<ul style="list-style-type: none"> <li>Full, multi-rows</li> <li>Multi-master/local replica, lazy</li> </ul>	<ul style="list-style-type: none"> <li>SER</li> <li>SI</li> <li>Session</li> </ul>	LWW	RO WO	Cost-consistency	<ul style="list-style-type: none"> <li>provides SSER and SI using 2PC, <i>group communication</i>, and logical clock.</li> <li>is not fault-tolerant.</li> </ul>
<b>Gemini</b> [Li et al. 2012] MDCs	<ul style="list-style-type: none"> <li>Full, multi-shard</li> <li>Multi-master, lazy/eager</li> </ul>	<ul style="list-style-type: none"> <li>Eventual</li> <li>Strong</li> </ul>	n/a	RO WO	Availability-consistency	<ul style="list-style-type: none"> <li>partitions operations based on the system semantic so that increases commutativity operations.</li> <li>does not tolerate faults.</li> </ul>
<b>Pileus</b> [Terry et al. 2013] MDCs	<ul style="list-style-type: none"> <li>Full, single row</li> <li>Single master, lazy</li> </ul>	<ul style="list-style-type: none"> <li>Weak**</li> <li>Strong</li> <li>RYW, MR</li> </ul>	n/a±	RW	Consistency-latency	<ul style="list-style-type: none"> <li>supports strong consistency via directing all reads to master replica and using a logical clock.</li> <li>is not fault-tolerant.</li> </ul>
<b>Bismar</b> [Chihoub et al. 2013] MDCs	<ul style="list-style-type: none"> <li>Full, single row</li> <li>Multi-master, lazy</li> </ul>	<ul style="list-style-type: none"> <li>Weak</li> <li>Quorum</li> <li>Strong</li> </ul>	n/a	n/a	Cost-latency	<ul style="list-style-type: none"> <li>demonstrates as consistency level is high cost and rate of fresh reads is high.</li> <li>is not tolerant to faults.</li> </ul>
<b>QUELEA</b> [Sivaramakrishnan et al. 2015] MDCs	<ul style="list-style-type: none"> <li>Full, multi-row</li> <li>Multi-master, lazy/eager</li> </ul>	<ul style="list-style-type: none"> <li>Weak ◊</li> <li>RC, RR, MAV◊◊</li> </ul>	LWW	n/a	Availability-consistency	<ul style="list-style-type: none"> <li>is a declarative language that provides different levels of consistency depending to operations conducted on top of eventually consistent data store.</li> </ul>
<b>Zhao et al. [2015]</b> SDC/MDCs	<ul style="list-style-type: none"> <li>Full, multishard</li> <li>Single-master, lazy</li> </ul>	<ul style="list-style-type: none"> <li>Strong</li> <li>Weak</li> </ul>	n/a	n/a	Consistency-latency Consistency-throughput	<ul style="list-style-type: none"> <li>propose dynamic provisioning mechanisms based on the SLA requirements, namely data freshness and requests response times, for Database-as-a-Services.</li> </ul>

† Abbreviations: Tan. (Transaction), Stict serializability (SSER), Snapshot Isolation (SI), Monotonic Read (MR), Read-Your-Write (RYW), Monotonic Atomic View (MAV), and Last Write Win (LWW).

†† Read-Your-Write (RYW) is also called *session* consistency. \* Similar to Sinfonia [Aguilera et al. 2007], CRAQ supports mini-transactions defined as a compare, read, and write set.

\*\* Weak consistency includes eventual and causal consistency modes. ± Pileus avoids conflict through using primary replica. ◊ Monotonic Atomic View (MAV) is defined as: if some operations in the transaction  $T_i$  observes the effects of another transaction  $T_j$ , then subsequent operations of  $T_i$  will observe the effects of  $T_j$ .

◊◊ Weak consistency model consists of all its models: eventual, causal, and ordering.



Table V: Summary of Projects with Strong Isolation Level **across DCs** in several aspects: Project specifications (name and data model), Transaction support (architecture, granularity, and type), Transaction protocol, Replication protocol, Contemporary technique, Commit latency (Message at bottleneck replica), main feature and drawback.

Project specifications	Transaction support	Isolation level	Transaction protocol	Rep.* protocol	Cont. technique	CL (MAB)	Feature	Drawback
<b>Spanner</b> <sup>†</sup> [Corbett et al. 2013] Semi-RDB	<ul style="list-style-type: none"> <li>• Top-down</li> <li>• Multi-shard</li> <li>• RO and RW<sup>††</sup></li> </ul>	Strict SER	2PL+2PC	Multi-Paxos	×	4 (2n)	RO transaction Optimization	Low throughput
<b>TARIP</b> [Zhang et al. 2015] KV	<ul style="list-style-type: none"> <li>• Top-down</li> <li>• Multi-shard</li> <li>• RO and RW</li> </ul>	Strict SER	Timestamp ordering +OCC	Inconsistent Replication (IR)	×	2 (2)	Latency optimization	Unsuitable for workloads with high contention
<b>Janus</b> [Mu et al. 2016] KV	<ul style="list-style-type: none"> <li>• Top-down</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	Strict SER	Dependency-tracking	Gen. Paxos	✓	2/4 (2)	Throughput optimization	Execution of <i>one-shot</i> transactions
<b>Clock-RSM</b> [Du et al. 2014b] KV	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	Strict SER	Physical timestamp ordering	SMR (all leader)	×	2 ( $n^2$ )	Latency optimization	Low throughput
<b>MegaStore</b> [Baker et al. 2011] Column families	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Single shard</li> <li>• RW</li> </ul>	SER	Paxos <sup>‡‡</sup> +2PC	Multi-Paxos	×	4 (2n)	Pioneer in guaranteeing transactions across shards	Low throughput
<b>Calvin</b> [Thomson et al. 2012] KV	<ul style="list-style-type: none"> <li>• Reversed top-down</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	SER	Serialization of transactions +(resemblance to)2PC	Paxos	×	4 (2n)	Latency optimization	Low throughput
<b>Paxos-CP</b> [Patterson et al. 2012] MDCs KV	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Single shard</li> <li>• RW</li> </ul>	SER	OCC+Paxos	Multi-Paxos	×	2 (2n)	Throughput improvement compared to Clock-RSM, MegaStore, and Calvin	Still suffering from low throughput
<b>Replicated Commit</b> [Mahmoud et al. 2013] KV	<ul style="list-style-type: none"> <li>• Reversed top-down</li> <li>• Multi-shard</li> <li>• RO and RW</li> </ul>	SER	2PL+2PC	Paxos	×	4 (2)	Reduction in RTT across DCs	High latency (2n) for RO transactions

\*Abbreviations: Rep. (Replication), CL (Commit Latency), and MAB (Message at Bottleneck).

† F1 [Shute et al. 2013] uses the same concurrency mechanism on the top of Spanner to support centralized and distributed SQL queries for OLTP and OLAP applications.

†† Spanner supports *snapshot reads* that execute a read operation without locking and retrieve the data based the timestamp specified by client.

‡ One-shot transactions do not contain any piece whose input is dependant on the execution of another piece.

‡‡ Paxos is used as concurrency control via the replicated log accessed serially by transactions.

Table V: Projects with Strong Isolation Level across DCs

Project specifications	Transaction support	Isolation level	Transaction protocol	Replication protocol	Cont. technique	CL (MAB)	Feature	Drawback
<b>Message Futures</b> [Nawab et al. 2013] KV	<ul style="list-style-type: none"> <li>Flat</li> <li>Multi-shard</li> <li>RW</li> </ul>	SER	Replicated log + serialization graphs [Bernstein et al. 1986]	Replicated Log across all DCs	×	2 (2n)	High throughput compared to Paxos-CP and Paxos	A high message transmission overhead across DCs
<b>CloudTPS</b> [Wei et al. 2012] KV	<ul style="list-style-type: none"> <li>Top-down</li> <li>Single shard</li> <li>RO and RW</li> </ul>	SER	Timestamp-ordering+2PC	Full Eager	×	5 (n)	Adaptable to SimpleDB and Bigtable	Low scalability and high memory overheads
<b>Deuteronomy</b> [Levandowski et al. 2011] KV	<ul style="list-style-type: none"> <li>Flat</li> <li>Multi-shard</li> <li>RW</li> </ul>	SER	locking-based ordering-disk-based log	Full Eager	×	4 (2n)	Separation of transaction manager and data store	Not scalable due to using a centralized transaction manager
<b>Lynx</b> [Zhang et al. 2013] RDB	<ul style="list-style-type: none"> <li>Flat</li> <li>Single shard</li> <li>RW</li> </ul>	SER	Transaction decomposition along with (2PL+2PC) and commutativity	Chain replication	✓	2n (2)	Scalable with the number of both servers and DCs	Requirement for a prior knowledge of transactions and static analysis
<b>P-store</b> [Schiper et al. 2010] KV	<ul style="list-style-type: none"> <li>Flat</li> <li>Multi-shard</li> <li>RW</li> </ul>	SER	Genuine atomic multicast+ 2PC-like protocol	Genuine partial replication (GPR)	×	4/5 (2n)	Providing the advantages in Table IX	Requirement for an expensive certification phase for RO transactions
<b>Clock-SI</b> [Du et al. 2013a] KV	<ul style="list-style-type: none"> <li>Flat</li> <li>Multi-shard</li> <li>RO and RW</li> </ul>	SI	Physical timestamp ordering +2PC	Full Eager	×	2/4 (2n)	Low transaction latency and high throughput††	Snapshot unavailability due to physical clock skew
<b>Walter</b> [Sovran et al. 2011] KV	<ul style="list-style-type: none"> <li>Flat</li> <li>Multi-shard</li> <li>RW</li> </ul>	PSI	Timestamp ordering +2PC	Full Lazy	✓	0/2 (2n)◇	Providing strong consistency within a DC and causal consistency across DCs	Requirement for a certification phase for RO transactions
<b>Jessy</b> [Ardekani et al. 2013] KV	<ul style="list-style-type: none"> <li>Flat</li> <li>Multi-shard</li> <li>RO and RW</li> </ul>	NMSI	Partial vector dependency (PVD)+genuine atomic multicast	GPR	×	4/5 (2n)	Comparable with Read committed (RC) isolation in latency and throughput of transactions	See Table XV
<b>MDCC</b> [Kraska et al. 2013] KV	<ul style="list-style-type: none"> <li>Flat</li> <li>Multi-shard</li> <li>RW</li> </ul>	RC	n/a	Gen. Paxos	✓	3 (2n)	Achievement of strong consistency at a cost similar to eventual consistency	See Table XV

† A transaction is global if its execution is distributed at different geographical locations.

†† Clock-SI achieve these features compared to the conventional Snapshot Isolation (i.e., referring to the implementation of snapshot isolation using a centralized timestamp authority) by reduction in 1 RTT for RO transactions and 2 RTTs for RW transactions.

◇ n is the number of DCs, not the number of replicas of the object.

Table VI: Summary of Projects with Strong Isolation Level **Within DCs** in several aspects: Project specifications, Transaction support (architecture, granularity, and type), Isolation level, Transaction protocol, Replication protocol, Replication degree, main feature and drawback.

Project specifications*	Transaction support	Isolation level	Transaction protocol	Rep. protocol	Rep. degree	Feature	Drawback
<b>Sinfonia</b> [Aguilera et al. 2007]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Single shard</li> <li>• RO and RW</li> </ul>	SER	OCC+(modified)2PC	Single master Lazy	Full	Providing mini-transactions	Requirement for a prior knowledge of reads and writes in transactions
<b>G-Store</b> [Das et al. 2010]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Single shard</li> <li>• RW</li> </ul>	SER (ACID)	OCC+write-ahead logging	n/a	n/a	Providing serializable transactions on top of HBase within a server	Not optimization for guarantee transactions across DCs
<b>Granola</b> [Cowling and Liskov 2012]	<ul style="list-style-type: none"> <li>• Top-down</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	SER	Timestamp/locking-based ordering	Viewstamped [Oki and Liskov 1988]	Full	High throughput and low latency for <i>independent transactions</i>	Requirement for a prior knowledge of reads and writes in transactions
<b>S-DUR</b> [Sciascia et al. 2012]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	SER	OCC+ atomic broadcast	DUR	Partial	The scalability of RO and local RW transactions with replicas number	High message exchange rate
<b>SCORE</b> [Peluso et al. 2012a]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Multi-shard</li> <li>• RO and RW</li> </ul>	SER	MVCC+genuine atomic	GPR	Partial	Providing abort-free RO transaction without a need for distributed certification phase	Requirement for a garbage collection mechanism to deal with obsolete data
<b>P-DUR</b> [Pacheco et al. 2014]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Multi shards</li> <li>• RW</li> </ul>	SER	Atomic multicast+2PC-like	DUR	Partial	The scalability of RW transactions with the number of cores available in a replica	Not scalable with the number of replicas
<b>MaaT</b> [Mahmoud et al. 2014]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	SER	A re-designed OCC++	Multi-master Eager	Full	Removing the drawbacks of existing OCC mechanism, while preserving the benefits of OCC over lock-based concurrency control (see section 6.1.2)	A source of overhead of memory requirement for read/write timestamp, (soft) locks table, and time table in each server.
<b>Warp</b> [Escriva et al. 2015]	<ul style="list-style-type: none"> <li>• Top-down</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	SER	Dependency-tracking	Chain replication [van Renesse and Schneider 2004]	Full	Execution of transactions in natural their arrival order unless doing so would violate SER.	Not optimization for Geo-replication.
<b>ElasTras</b> [Das et al. 2013] KV	<ul style="list-style-type: none"> <li>• Top-down</li> <li>• Single shard</li> <li>• RW</li> </ul>	SER	OCC+ a log in each Transaction manager	Multi-master Eager	Full	Improving multi-tenancy and elasticity	Providing a limited distributed transactions

\* All projects use *key-value* data model, and they neither deploy contemporary technique to guarantee isolation level.

† It uses *dynamic timestamp* to eliminate distributed *certification procedure*, and *soft locks*, to inform transactions accessing a data item which is read or written by other uncommitted transactions.

Table VI: Projects with Strong Isolation Level **within DCs**

Project specifications	Transaction support	Isolation level	Transaction protocol	Rep. protocol	Rep. degree	Features	Drawbacks
<b>GMU</b> [Peluso et al. 2012b]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	EUS	MVCC+2PC	GPR	Partial	Neither a need for centralized and global clock nor requirement for any remote certification phase for RO transactions	Incurring the occurrence of distributed deadlocks, which results in poor scalability
<b>Percolator</b> [Peng and Dabek 2010]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	SI	MVCC+2PC	Multi-Paxos <sup>†</sup>	n/a	Aiming at the batched execution model where low transaction latency is not a goal	Uses a central server to issue timestamps to transactions
<b>Omid</b> [Gomez Ferro et al. 2014]	<ul style="list-style-type: none"> <li>• Flat</li> <li>• Multi-shard</li> <li>• RW</li> </ul>	SI	Timestamp ordering	n/a	Full	Unlike Percolator, avoiding multi-version overheads and deadlock	Not scalable due to using a centralized timestamp-ordering

<sup>†</sup> Percolator was built on top of Bigtable which uses a Multi-paxos replication protocol.

## F. PROJECT WITH MONETARY COST OPTIMIZATION

Table VII: Summary of Projects with Monetary Cost Optimization in Four Aspects: (i) Project specification (name, data store platform, data application type), (ii) cost elements, (iii) QoS metrics/constraints, (iv) Features (solution techniques, cost optimization type, and key feature).

Project specification	Cost elements					QoS metrics/constraints					Features
	Storage	Read	Write	VM	Migration	Latency	Bandwidth	Availability	durability	Consistency level	
[Chang et al. 2012] MDC Data-intensive	✓							✓			(i) Dynamic programming. (ii) Cost-based availability optimization. (iii) Maximizing availability under budget constraint for non-partitioned objects
[Mansouri et al. 2013] MDC Data-intensive	✓							✓			(i) Dynamic programming and Brute-force search. (ii) Cost-based availability optimization. (iii) Minimizing cost under a given availability for non-partitioned objects and maximizing availability under a given budget for partitioned objects.
PRCR [Li et al. 2012] SDC Scientific workflows	✓								✓		(i) Proactive replica checking. (ii) Cost-based durability optimization. (iii) Reducing $\frac{1}{3}$ to $\frac{2}{3}$ storage cost as compared triplicate storage cost.
CIR [Li et al. 2011] SDC Scientific workflows	✓								✓		(i) an incremental replication approach. (ii) Cost-based durability optimization. (iii) Incurring cost the same as that for triplicate replicas in long-term storage deployment.
Copysets [Cidon et al. 2013] SDC Data-intensive	✓								✓		(i) Near optimal solution. (ii) Cost-based durability optimization. (iii) Achieving higher durability than widely used random replication
Consistency rationing [Kraska et al. 2009] SDC OLTP	✓									✓	(i) A general policy based on a conflict probability (ii) Cost-based consistency optimization. (iii) Providing <i>serializability</i> , <i>session</i> , and <i>adaptive consistency</i> between these two consistency models.
[Mansouri and Buyya 2016] a dual DC Data-intensive	✓	✓	✓		✓	✓					(i) Dynamic programming and greedy algorithms. (ii) Storage-network trade-off. (iii) Providing a optimization cost for objects in OSN.
C3 [Fetai and Schuldt 2012] SDC Data-intensive	✓	✓	✓	✓						✓	(i) A selective solution based on defined rules. (ii) Cost-based consistency optimization. (iii) Providing <i>serializability</i> , <i>snapshot isolation</i> , and <i>eventual consistency</i> .
[Narayanan et al. 2014] MDC OSN	✓	✓	✓			✓				✓	(i) Linear programming (ii) Cost-based latency optimization. (iii) Minimizing latency in quorum-based data stores and optimizing monetary cost under latency constraint.
CosTLO [Wu et al. 2015b] MDC Data-intensive	✓	✓	✓	✓		✓				✓	(i) A comprehensive measurement study on S3 and Azure. (ii) Cost-based latency optimization. (iii) Reduction in latency variation via augmenting read/write requests with a 15% increase in cost.
[Puttaswamy et al. 2012] SDC File system	✓	✓	✓								(i) Ski-rental problem. (ii) Storage-cache trade-off. (iii) Deploying EBS and S3 to reduce monetary cost for a file system.
[Khanafer et al. 2013] SDC File system	✓	✓	✓								(i) A variant of ski-rental problem. (ii) Storage-cache trade-off. (iii) Reducing the deploying cost of file system across S3 and EBS as assumed the average and variance of arrival time of reads/ writes are known.
[Jiao et al. 2014] MDC OSN	✓	✓	✓			✓					(i) Using <i>graph cut</i> techniques to determine master and slave replicas. (ii) Cost optimization based on single QoS metric. (iii) Optimizing the mentioned costs plus carbon footprint cost.
Cosplay [Jiao et al. 2016] MDC OSN	✓		✓			✓		✓		Eventual	(i) A selective solution (greedy). (ii) Cost optimization based on multi-QoS metric. (iii) Achieving cost reduction via role-swaps between master and slave replicas in a greedy approach.
[Hu et al. 2016] MDC CDN	✓	✓	✓			✓					(i) Lynapunov optimization technique. (ii) Cost optimization based on single QoS metric. (iii) Seeking the optimal solution without the future knowledge in regard to workload.

Table VII: Summary of Projects with Monetary Cost Optimization

Project specification	Cost elements					QoS metrics/constraints					Features
	Storage	Read	Write	VM	Migration	Latency	Bandwidth	Availability	durability	Consistency level	
[Chen et al. 2012] MDC CDN	✓	✓	✓			✓					(i) A selective solution (greedy). (ii) Cost optimization based on single QoS metric. (iii) Achieving cost reduction via greedy algorithms in the case of offline and online (Dynamic and Static) replication.
[Papagianni et al. 2013] MDC CDN	✓	✓	✓	✓		✓	✓				(i) Linear programming and graph partitioning heuristic algorithms. (ii) Cost optimization based on multi-QoS metric. (iii) Reducing cost by considering the optimized replica placement and distribution path construction.
COMIC [Yao et al. 2015] MDC CDN	✓	✓									(i) Mixed integer linear programming. (ii) Cost optimization based on Multi-QoS metric. (iii) Optimizing electricity cost for DCs and usage CDN cost (read cost) under the processing capacity of DCs and CDN as constraints.
[Wu et al. 2015a] MDC social media streaming	✓	✓		✓	✓	✓	✓				(i) An integer programming and an online algorithm based on prediction. (ii) Cost optimization based on multi-QoS metric. (iii) Reducing cost per video as workload changes across DCs.
[Qiu et al. 2015] MDC (public and private) social media streaming	✓	✓		✓	✓	✓	✓				(i) Lynapunov optimization technique. (ii) Cost optimization based on multi-QoS metric. (iii) Reducing cost (migration cost is only between private and public DC) for video files while ensuring bandwidth constraint for the private DC.
[Ruiz-Alvarez and Humphrey 2012] MDC (Private and Amazon) Data-intensive	✓	✓	✓	✓		✓	✓				(i) Integer linear programming. (ii) Cost optimization based on multi-QoS metric. (iii) Determining either private or Amazon cloud to run application based on budget and job turnaround time besides of the mentioned constraints.
SPANStore [Wu et al. 2013] MDC Data-intensive	✓	✓	✓	✓		✓		✓		Eventual Strong	(i) Linear programming. (ii) Cost optimization based on multi-QoS metric. (iii) VM cost optimization only for writes propagation.
[Chiu and Agrawal 2010] SDC Data-intensive	✓	✓	✓	✓							(i) A comprehensive scenario conducted on S3, EBS and EC2 instances. (ii) Storage-cache trade-off. (iii) A reduction cost via using S3 vs. EBS and cache owning by different EC2 instance.
[Yuan et al. 2011; Yuan et al. 2013] SDC Scientific workflows	✓			✓							(i) Using intermediate data dependency graph or Cost Transitive Tournament Shortest Path (CTT-SP)-based to decide whether to store data or recompute later in run time. (ii) Storage-computation trade-off. (iii) Significant reduction in cost by using Amazon cost model.
[Jokhio et al. 2013] SDC Video Transcoding	✓			✓							(i) A selective solution utilizing estimation of storage and computation costs and the popularity of transcoded video. (ii) Storage-computation trade-off. (iii) Significant reduction in cost by using Amazon cost model.
[Byholm et al. 2015] SDC Video Transcoding	✓			✓							(i) Using utility-based model which determines when and for how long each transcoded video should be stored. (ii) Storage-computation trade-off. (iii) Decision of model is based on the storage duration $t$ , the average of arrival request during $t$ , and the popularity distribution of video.
[Deelman et al. 2008] SDC data-intensive	✓	✓	✓	✓							(i) Just measuring the incurred cost for data-intensive application. (ii) Storage-computation trade-off. (iii) storing data for long term is more cost-effective than recomputing it later (about next two years- for this specific application).
Triones [Su et al. 2016] MDC data-intensive	✓	✓	✓		✓						(i) Non-linear programming and geometric space. (ii) Storage-computation trade-off. (iii) Improving fault-tolerance (2 times) and reducing access latency and vendor lock-in at the expense of more monetary cost.

## REFERENCES

- Marcos K. Aguilera, Arif Merchant, Mehul Shah, Alistair Veitch, and Christos Karamanolis. 2007. Sinfonia: A New Paradigm for Building Scalable Distributed Systems. (2007).
- Sérgio Almeida, João Leitão, and Luís Rodrigues. 2013. ChainReaction: A Causal+ Consistent Datastore Based on Chain Replication (*EuroSys '13*).
- Masoud Saeida Ardekani, Pierre Sutra, and Marc Shapiro. 2013. Non-monotonic Snapshot Isolation: Scalable and Strong Consistency for Geo-replicated Transactional Systems (*SRDS '13*).
- Peter Bailis, Ali Ghodsi, Joseph M. Hellerstein, and Ion Stoica. 2013. Bolt-on Causal Consistency (*SIGMOD '13*).
- Jason Baker, Chris Bond, James C. Corbett, and et al. 2011. Megastore: Providing Scalable, Highly Available Storage for Interactive Services (*CIDR '11*).
- Nalini Belaramani, Mike Dahlin, Lei Gao, and et al. 2006. PRACTI Replication (*NSDI '06*).
- David Bermbach, Markus Klems, Stefan Tai, and Michael Menzel. 2011. MetaStorage: A Federated Cloud Storage System to Manage Consistency-Latency Tradeoffs. In *IEEE Conference on Cloud Computing (CLOUD)*, 2011.
- Philip A Bernstein, Vassos Hadzilacos, and Nathan Goodman. 1986. *Concurrency Control and Recovery in Database Systems*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA.
- Matthias Brantner, Daniela Florescu, David Graf, Donald Kossmann, and Tim Kraska. 2008. Building a Database on S3 (*SIGMOD '08*).
- Benjamin Byholm, Fareed Jokhio, Adnan Ashraf, and et al. 2015. Cost-Efficient, Utility-Based Caching of Expensive Computations in the Cloud. In *Euromicro Conference on Parallel, Distributed, and Network-Based Processing*, 2015.
- Chia-Wei Chang, Pangfeng Liu, and Jan-Jan Wu. 2012. Probability-Based Cloud Storage Providers Selection Algorithms with Maximum Availability. In *Proceedings International Conference on Parallel Processing (ICPP)*, 2012.
- Fangfei Chen, Katherine Guo, John Lin, and Thomas La Porta. 2012. Intra-cloud lightning: Building CDNs in the cloud. In *INFOCOM*, 2012.
- Houssein-Eddine Chihoub, Shadi Ibrahim, abriel Antoniu, and Maria Perez. 2013. Consistency in the Cloud: When Money Does Matter!. In *IEEE/ACM Symposium on Cluster, Cloud and Grid Computing (CCGrid)*, 2013.
- David Chiu and Gagan Agrawal. 2010. Evaluating caching and storage options on the Amazon Web Services Cloud. In *IEEE/ACM International Conference on Grid Computing*, 2010.
- Asaf Cidon, Stephen M. Rumble, Ryan Stutsman, and et al. 2013. Copysets: Reducing the Frequency of Data Loss in Cloud Storage (*USENIX ATC '13*).
- James C. Corbett, Jeffrey Dean, Michael Epstein, and et al. 2013. Spanner: Google's Globally Distributed Database. *ACM Transactions on Computer Systems* (2013).
- James Cowling and Barbara Liskov. 2012. Granola: Low-overhead Distributed Transaction Coordination (*USENIX ATC '12*).
- Sudipto Das, Divyakant Agrawal, and Amr El Abbadi. 2010. G-Store: A Scalable Data Store for Transactional Multi Key Access in the Cloud (*SoCC '10*).
- Sudipto Das, Divyakant Agrawal, and Amr El Abbadi. 2013. ElasTraS: An Elastic, Scalable, and Self-managing Transactional Database for the Cloud. *ACM Transactions on Database Systems* (2013).
- E. Deelman, G. Singh, M. Livny, B. Berriman, and J. Good. 2008. The cost of doing science on the cloud: The Montage example. In *Conference for High Performance Computing, Networking, Storage and Analysis*, 2008.
- Jiaqing Du, Sameh Elnikety, Amitabha Roy, and Willy Zwaenepoel. 2013b. Orbe: Scalable Causal Consistency Using Dependency Matrices and Physical Clocks (*SOCC '13*).
- Jiaqing Du, S. Elnikety, and W. Zwaenepoel. 2013a. Clock-SI: Snapshot Isolation for Partitioned Data Stores Using Loosely Synchronized Clocks (*SRDS '13*).
- Jiaqing Du, Călin Iorgulescu, Amitabha Roy, and Willy Zwaenepoel. 2014a. Gentlerain: Cheap and scalable causal consistency with physical clocks. In *Proceedings of the ACM Symposium on Cloud Computing*, 2014.
- Jiaqing Du, Daniele Sciascia, Sameh Elnikety, Willy Zwaenepoel, and Fernando Pedone. 2014b. Clock-RSM: Low-Latency Inter-datacenter State Machine Replication Using Loosely Synchronized Physical Clocks. In *IEEE/IFIP Conference on Dependable Systems and Networks (DSN)*, 2014.
- Robert Escriva, Bernard Wong, and Emin Gün Sirer. 2015. Warp: Multi-key transactions for keyvalue stores. (2015).
- Ilir Fetai and Heiko Schuldt. 2012. Cost-Based Data Consistency in a Data-as-a-Service Cloud Environment. In *IEEE Conference on Cloud Computing (Cloud)*, 2012.
- Daniel Gomez Ferro, Flavio Junqueira, Ivan Kelly, Benjamin Reed, and Maysam Yabandeh. 2014. Omid: Lock-free transactional support for distributed data stores. In *IEEE International Conference on Data Engineering*, 2014.
- Han Hu, Yonggang Wen, Tat-Seng Chua, and et al. 2016. Joint Content Replication and Request Routing for Social Video Distribution over Cloud CDN: A Community Clustering Method. *IEEE Transactions on Circuits and Systems for Video Technology* (2016).
- Lei Jiao, Jun Li, Tianyin Xu, Wei Du, and Xiaoming Fu. 2016. Optimizing Cost for Online Social Networks on Geo-Distributed Clouds. *IEEE/ACM Transactions on Networking* (2016).
- Lei Jiao, Jun Lit, Wei Du, and Xiaoming Fu. 2014. Multi-objective data placement for multi-cloud socially aware services. In *INFOCOM*, 2014.
- F. Jokhio, A. Ashraf, S. Lafond, and J. Lilius. 2013. A Computation and Storage Trade-off Strategy for Cost-Efficient Video Transcoding in the Cloud. In *Euromicro Conference on Software Engineering and Advanced Applications*, 2013.
- Ali Khanafer, Murali Kodialam, and Krishna P. N. Puttaswamy. 2013. The constrained Ski-Rental problem and its application to online cloud cost optimization. In *INFOCOM*, 2013.
- Tim Kraska, Martin Hentschel, Gustavo Alonso, and Donald Kossmann. 2009. Consistency Rationing in the Cloud: Pay Only when It Matters. *Proceedings of the VLDB Endowment* (2009).



- Tim Kraska, Gene Pang, Michael J. Franklin, Samuel Madden, and Alan Fekete. 2013. MDCC: Multi-data Center Consistency (*EuroSys '13*).
- Justin J. Levandoski, David B. Lomet, Mohamed F. Mokbel, and Kevin Zhao. 2011. Deuteronomy: Transaction Support for Cloud Data (*CIDR '11*).
- Cheng Li, Daniel Porto, Allen Clement, and et al. 2012. Making Geo-replicated Systems Fast As Possible, Consistent when Necessary (*OSDI '12*).
- Wenhao Li, Yun Yang, Jinjun Chen, and Dong Yuan. 2012. A Cost-Effective Mechanism for Cloud Data Reliability Management Based on Proactive Replica Checking. In *Proceedings IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID), 2012*.
- Wenhao Li, Yun Yang, and Dong Yuan. 2011. A Novel Cost-Effective Dynamic Data Replication Strategy for Reliability in Cloud Data Centres. In *IEEE Conference on Dependable, Autonomic and Secure Computing (DASC), 2011*.
- Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Andersen. 2011. Don'T Settle for Eventual: Scalable Causal Consistency for Wide-area Storage with COPS (*SOSP '11*).
- Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Andersen. 2013. Stronger Semantics for Low-latency Geo-replicated Storage (*NSDI '13*).
- Hatem Mahmoud, Faisal Nawab, Alexander Pucher, Divyakant Agrawal, and Amr El Abbadi. 2013. Low-latency Multi-datacenter Databases Using Replicated Commit. *Proceedings of the VLDB Endowment* (2013).
- Hatem A. Mahmoud, Vaibhav Arora, Faisal Nawab, Divyakant Agrawal, and Amr El Abbadi. 2014. MaaT: Effective and Scalable Coordination of Distributed Transactions in the Cloud. *Proceedings of the VLDB Endowment* (2014).
- Yaser Mansouri and Rajkumar Buyya. 2016. To Move or Not to Move: Cost Optimization in a Dual Cloud-based Storage Architecture. *Journal of Network and Computer Applications* (2016).
- Yaser Mansouri, Adel Nadjaran Toosi, and Rajkumar Buyya. 2013. Brokering Algorithms for Optimizing the Availability and Cost of Cloud Storage Services. In *Conference on Cloud Computing Technology and Science (Cloudcom), 2013*.
- Shuai Mu, Lamont Nelson, Wyatt Lloyd, and Jinyang Li. 2016. Consolidating Concurrency Control and Consensus for Commits under Conflicts. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. USENIX Association, GA, 517–532.
- Shankaranarayanan Puzhavakath Narayanan, Ashiwan Sivakumar, Sanjay Rao, and Mohit Tawarmalani. 2014. Performance Sensitive Replication in Geo-distributed Cloud Datastores. In *IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), 2014*.
- Faisal Nawab, Divyakant Agrawal, and Amr El Abbadi. 2013. Message Futures: Fast Commitment of Transactions in Multi-datacenter Environments. In *Sixth Biennial Conference on Innovative Data Systems Research (CIDR), 2013*.
- Faisal Nawab, Vaibhav Arora, Divyakant Agrawal, and AE Abbadi. 2015. Chariots: A scalable shared log for data management in multi-datacenter cloud environments. In *Conference on Extending Database Technology (EDBT), 2015*.
- Brian M. Oki and Barbara H. Liskov. 1988. Viewstamped Replication: A New Primary Copy Method to Support Highly-Available Distributed Systems (*PODC '88*).
- Leandro Pacheco, Daniele Sciascia, and Fernando Pedone. 2014. Parallel Deferred Update Replication. In *IEEE Symposium on Network Computing and Applications (NCA), 2014*.
- Chrysa Papagianni, Aris Leivadeas, and Symeon Papavassiliou. 2013. A Cloud-Oriented Content Delivery Network Paradigm: Modeling and Assessment. *IEEE Transactions on Dependable and Secure Computing* (2013).
- Stacy Patterson, Aaron J. Elmore, Faisal Nawab, Divyakant Agrawal, and Amr El Abbadi. 2012. Serializability, Not Serial: Concurrency Control and Availability in Multi-datacenter Datastores. *Proceedings of the VLDB Endowment* (2012).
- Sebastiano Peluso, Paolo Romano, and Francesco Quaglia. 2012a. SCORE: A Scalable One-copy Serializable Partial Replication Protocol. In *Proceedings of the 13th International Middleware Conference (Middleware '12)*.
- Sebastiano Peluso, Pedro Ruivo, Paolo Romano, Francesco Quaglia, and Luis Rodrigues. 2012b. When Scalability Meets Consistency: Genuine Multiversion Update-Serializable Partial Data Replication. In *ICDCS, 2012*.
- Daniel Peng and Frank Dabek. 2010. Large-scale Incremental Processing Using Distributed Transactions and Notifications (*OSDI '10*).
- Krishna P.N. Puttaswamy, Thyaga Nandagopal, and Murali Kodialam. 2012. Frugal Storage for Cloud File Systems (*EuroSys '12*).
- Xuanjia Qiu, Hongxing Li, Chuan Wu, Zongpeng Li, and Francis C.M. Lau. 2015. Cost-Minimizing Dynamic Migration of Content Distribution Services into Hybrid Clouds. *IEEE Transactions on Parallel and Distributed Systems* (2015).
- Arkaitz Ruiz-Alvarez and Marty Humphrey. 2012. A Model and Decision Procedure for Data Storage in Cloud Computing. In *IEEE/ACM Symposium on Cluster, Cloud and Grid Computing (CCGRID), 2012*.
- Nicolas Schiper, Pierre Sutra, and Fernando Pedone. 2010. P-Store: Genuine Partial Replication in Wide Area Networks (*SRDS '10*).
- Daniele Sciascia, Fernando Pedone, and Flavio Junqueira. 2012. Scalable Deferred Update Replication (*DSN '12*).
- Min Shen, Ajay D. Kshemkalyani, and Ta yuan Hsu. 2015. Causal Consistency for Geo-Replicated Cloud Storage under Partial Replication. In *Parallel and Distributed Processing Symposium Workshop (IPDPSW), 2015*.
- Jeff Shute, Radek Vingralek, Bart Samwel, and et al. 2013. F1: A Distributed SQL Database That Scales. *Proceedings of the VLDB Endowment* (2013).
- KC Sivaramakrishnan, Gowtham Kaki, and Suresh Jagannathan. 2015. Declarative Programming over Eventually Consistent Data Stores (*PLDI '15*).
- Yair Sovran, Russell Power, Marcos K. Aguilera, and Jinyang Li. 2011. Transactional Storage for Geo-replicated Systems (*SOSP '11*).
- Maomeng Su, Lei Zhang, Yongwei Wu, Kang Chen, and Keqin Li. 2016. Systematic Data Placement Optimization in Multi-Cloud Storage for Complex Requirements. *IEEE Trans. Comput.* (2016).
- Jeff Terrace and Michael J. Freedman. 2009. Object Storage on CRAQ: High-throughput Chain Replication for Read-mostly Workloads (*USENIX '09*).

- Douglas B. Terry, Vijayan Prabhakaran, Ramakrishna Kotla, and et al. 2013. Consistency-based service level agreements for cloud storage (*SOSP '13*).
- Alexander Thomson, Thaddeus Diamond, Shu-Chun Weng, and et al. 2012. Calvin: Fast Distributed Transactions for Partitioned Database Systems (*SIGMOD '12*).
- Robbert van Renesse and Fred B. Schneider. 2004. Chain Replication for Supporting High Throughput and Availability (*OSDI '04*).
- Zhou Wei, G. Pierre, and Chi-Hung Chi. 2012. CloudTPS: Scalable Transactions for Web Applications in the Cloud. *IEEE Transactions on Services Computing* (2012).
- Yu Wu, Chuan Wu, Bo Li, and et al. 2015a. Scaling Social Media Applications into Geo-distributed Clouds. *IEEE/ACM Transactions on Networking* (2015).
- Zhe Wu, Michael Butkiewicz, Dorian Perkins, Ethan Katz-Bassett, and Harsha V. Madhyastha. 2013. SPANStore: Cost-effective Geo-replicated Storage Spanning Multiple Cloud Services (*SOSP '13*).
- Zhe Wu, Curtis Yu, and Harsha V. Madhyastha. 2015b. CosTLO: Cost-effective Redundancy for Lower Latency Variance on Cloud Storage Services (*NSDI '15*).
- Jianguo Yao, Haihang Zhou, Jianying Luo, Xue Liu, and Haibing Guan. 2015. COMIC: Cost Optimization for Internet Content Multihoming. *IEEE Transactions on Parallel and Distributed Systems* (2015).
- Dong Yuan, Yun Yang, Xiao Liu, and Jinjun Chen. 2011. On-demand minimum cost benchmarking for intermediate dataset storage in scientific cloud workflow systems. *J. Parallel and Distrib. Comput.* (2011).
- Dong Yuan, Yun Yang, Xiao Liu, and et al. 2013. A Highly Practical Approach toward Achieving Minimum Data Sets Storage Cost in the Cloud. *IEEE Transactions on Parallel and Distributed Systems* (2013).
- Irene Zhang, Naveen Kr. Sharma, Adriana Szekeres, Arvind Krishnamurthy, and Dan R. K. Ports. 2015. Building Consistent Transactions with Inconsistent Replication (*SOSP '15*).
- Yang Zhang, Russell Power, Siyuan Zhou, and et al. 2013. Transaction Chains: Achieving Serializability with Low Latency in Geo-distributed Storage Systems (*SOSP '13*).
- L. Zhao, S. Sakr, and A. Liu. 2015. A Framework for Consumer-Centric SLA Management of Cloud-Hosted Databases. *IEEE Transactions on Services Computing* (July 2015).