Contents lists available at ScienceDirect

# Future Generation Computer Systems

# Multi-site data distribution for disaster recovery—A planning framework

Shubhashis Sengupta *, K.M. Annervaz

*Accenture Technology Labs, Bangalore, India*

## HIGHLIGHTS

- We describe a fault-tolerant multi-cloud data backup scheme using erasure coding.
- The data is distributed using a plan driven by a multi-criteria optimization.
- The plan uses parameters like cost, replication level, recoverability objective etc.
- Both single customer and multiple customer cases are tackled.
- Simulation results for the plans and sensitivity analyses are discussed.

## ARTICLE INFO

## ABSTRACT

In this paper, we present DDP-DR: a Data Distribution Planner for Disaster Recovery. DDP-DR provides an optimal way of backing-up critical business data into data centers (DCs) across several Geographic locations. DDP-DR provides a plan for replication of backup data across potentially large number of data centers so that (i) the client data is recoverable in the event of catastrophic failure at one or more data centers (disaster recovery) and, (ii) the client data is replicated and distributed in an optimal way taking into consideration major business criteria such as cost of storage, protection level against site failures, and other business and operational parameters like recovery point objective (RPO), and recovery time objective (RTO). The planner uses Erasure Coding (EC) to divide and codify data chunks into fragments and distribute the fragments across DR sites or storage zones so that failure of one or more site / zone can be tolerated and data can be regenerated. We describe data distribution planning approaches for both single customer and multiple customer scenarios.

## 1. Introduction

In today's enterprise computing, data centers generate an overwhelming volume of data. Applications such as particle physics [1], storing of web pages and indexes [2], social networking applications, and engineering applications of pharmaceutical and semiconductor companies can easily generate petabytes of data over days and weeks. Disaster Recovery (DR) and Business Continuity planning (BCP) require that critical enterprise data is backed up periodically and kept in geographically separate and secure locations. In the event of operational disruption at the primary site, the operation can be resumed at an alternate site where the backed up data

and log files are shipped and applications/services can be instantiated again. Additionally, recent regulatory and compliance standards like HIPAA, SOX and GLBA mandate that all operational data is retained for a certain period of time and be made available for auditing. With the increasing volume of data and increasing emphasis on service availability and data retention, the technology and process of handling backup and recovery have come under renewed scrutiny.

### 1.1. Traditional backup methodology

Traditionally, the data backup and archival are done using magnetic tapes which are processed and transported to a remote location. However, such procedure is manual and cumbersome (therefore slow) and rapid data restoration and service resumption are often not possible. Recently, with the advent of cheap, improved storage and online disk backup technology, and advances in networking; online remote backup options have become attractive

* Corresponding author. Tel.: +91 9845185828.
*E-mail addresses:* shubhashis.sengupta@accenture.com (S. Sengupta), annervaz.k.m@accenture.com (K.M. Annervaz).
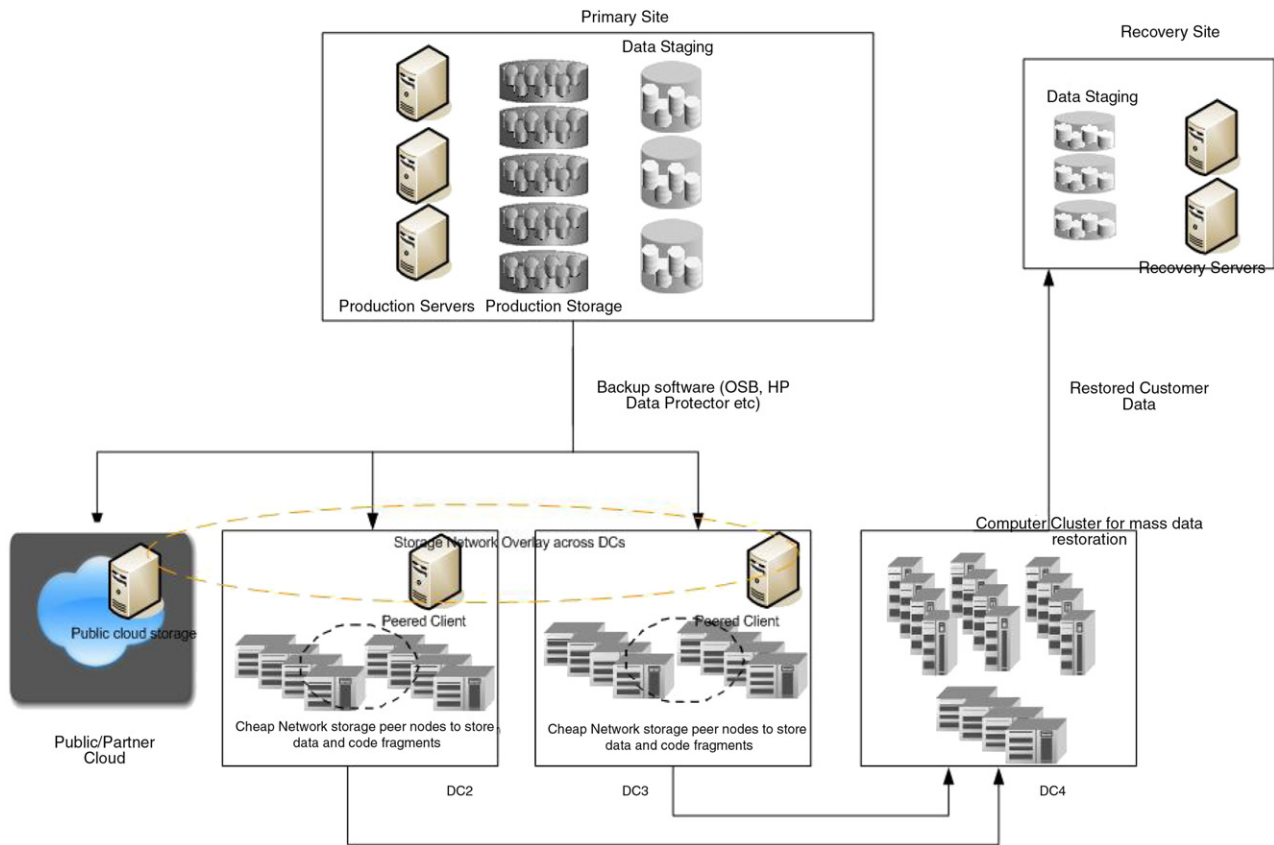
**Fig. 1.** Schematic multi-site DR infrastructure.

[3]. The storage area network and virtualization technology has become sophisticated enough to create a storage volume snapshot to a remote site [4]. Increasingly, open-source technologies such as RSync [5] are being used to achieve the same goals; albeit with a lower efficiency.

### 1.2. Multi-site data replication and backup

Cloud computing and cheap online storage technologies are promising to change the landscape of disaster recovery. The data from the primary site is now backed up in the cloud and/or in multiple geographically separated data centers to improve fault tolerance and availability [6]. Several cloud infrastructure and storage vendors such as Amazon S3, Glacier [7], and Rackspace [8] provide storage for backup. Several other vendors like Zamanda [9], use the cloud storage, such as Amazon S3, to provide backup services. Organizations are also adopting hybrid approach—where very critical or sensitive data is stored within the enterprise and non-sensitive data is dispatched to cloud. While backup using a single cloud or online storage is cheap and practical, online storing of encrypted backup data to a single third-party storage provider may not be prudent due to the lack of operational control, security, reliability and availability issues. It is advisable that organizations hedge their bets by replicating data to multiple cloud locations and data centers. It is also observed that in large organizations, having data centers (DC) in multiple geographies, DR may involve using one regional data center as an alternate site against another by replicating data. Replicating DR data across sites improves the availability by reducing the risk of simultaneous co-related failures.

In this context, we present a schematic diagram for a multi-site DR in Fig. 1. The primary site (DC1) hosts the servers and storage for production, test, and development. Historical operational data

is periodically copied to the staging servers where aggregation and de-duplication are run. The "backup" ready data is then replicated to multiple data centers[1] that the firm owns (DC2 and DC3) and/or to the public cloud storage providers. In the event of failure at the primary site, the data can be recovered to the recovery (also called secondary) site on demand. Data recovery or retrieval may require additional compute resources to carry out costly operations such as de-compression and decryption of data. Therefore, recovery can be optionally offloaded to a server firm (DC4) or to a dedicated processing hardware in the DR sites that can do bulk recovery of multiple customers within stipulated time bounds.

Since we propose a distributed storage substrate, one possible mechanism to maintain data consistency across backup sites is to create a peer-to-peer based storage overlay layer across the sites. Various distributed archival storage substrate are discussed in literature [10,11], but this is not the primary concern of this paper.

### 1.3. Optimal data distribution plan for multi-site backup (Motivation of this work)

The current approach of multi-site backup is to replicate data to single or multiple remote sites so that co-related storage or network failures do not hamper data availability. Replication, however, increases data redundancy linearly with the number of sites. Plain replication, even with data compression technologies, makes the data footprint quite large. Additionally, it is often seen that the strategy of data placement and distribution is not driven

---

[1] We use the term data center in a broad sense. It may also mean a set of storage nodes/cluster.