

# The Best Practices and Performance Improvements of Cloud Storage

## China Railway ITC



## BUSINESS SITUATION

### China Railway

A solely state-owned enterprise

The main artery of the national economy

#### MAIN BUSINESS

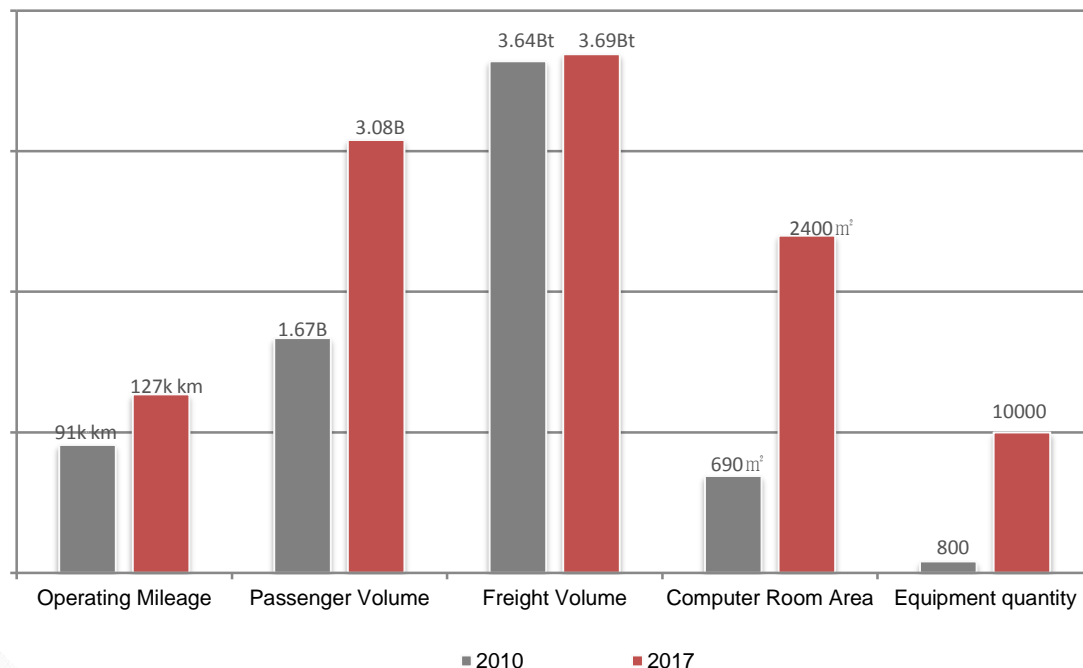
Passenger and freight transport service

#### BUSINESS FEATURE

Large scale, wide coverage, uninterrupted

#### ENTERPRISE GOAL

A world-class modern logistics enterprise



## Cloud Computing Development in China Railway

### Starting in 2014

Build the cloud data center in stages, and gradually migrated production applications to the cloud environment.



### OpenStack-base for Cloud Data Center

Currently reached a scale of thousands of physical machines, and finished deploying dozens of applications, including passenger transport, freight, scheduling, locomotive and public infrastructure platform.



### Cloud Computing Powering Data Center Hub

Expected to reach the scale of newly built Data Center Hub with above 15,000 physical machines by the end of 2018.

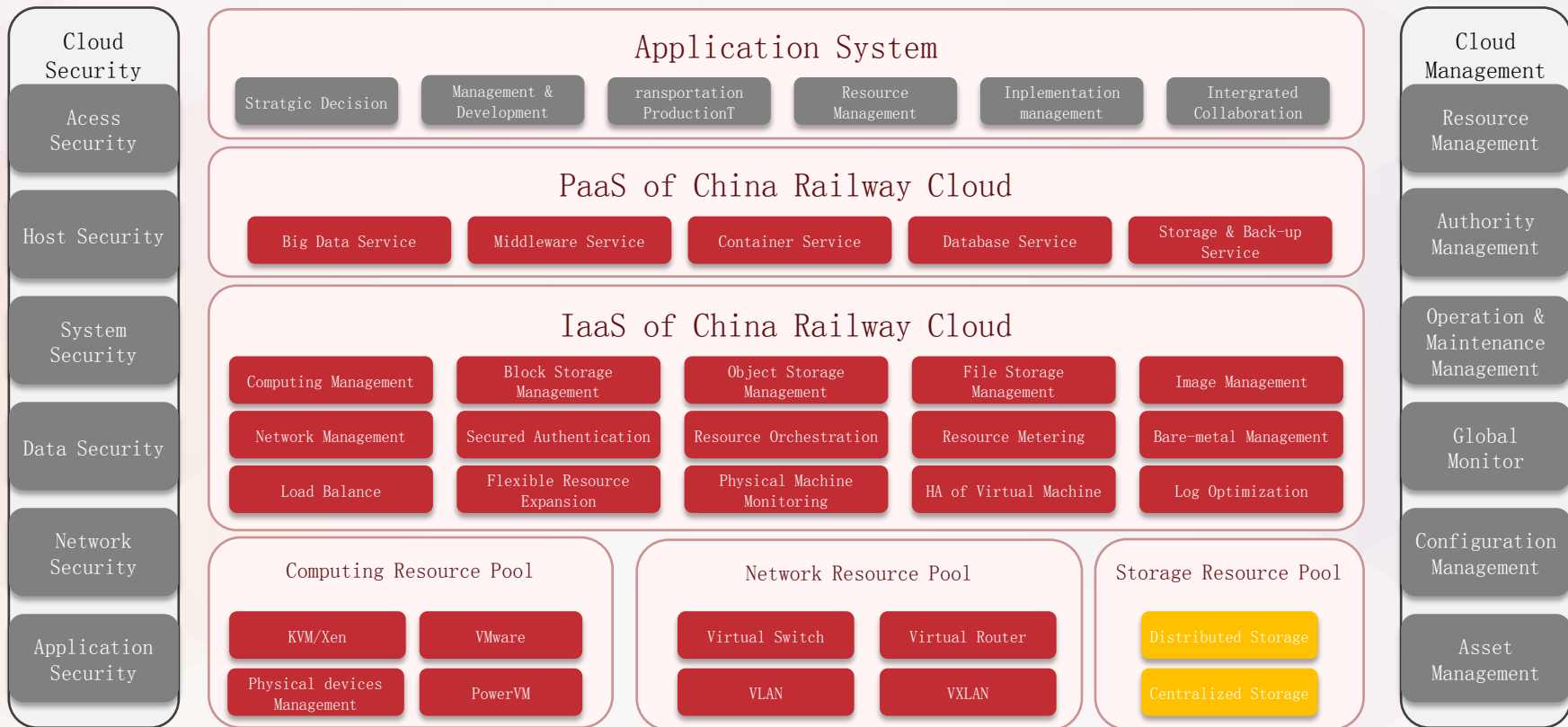


---

As shown in this figure, it is the current storage usage of China Railway data center, where the centralized storage services for the database service and the back-end storage of VMware. Both centralized storage and distributed storage serve as the back-end storage of the cloud platform. Due to our lack of familiarity with distributed storage, we also want to embrace open source and pursue an open approach, so we temporarily put some light applications on distributed storage, such as web server.

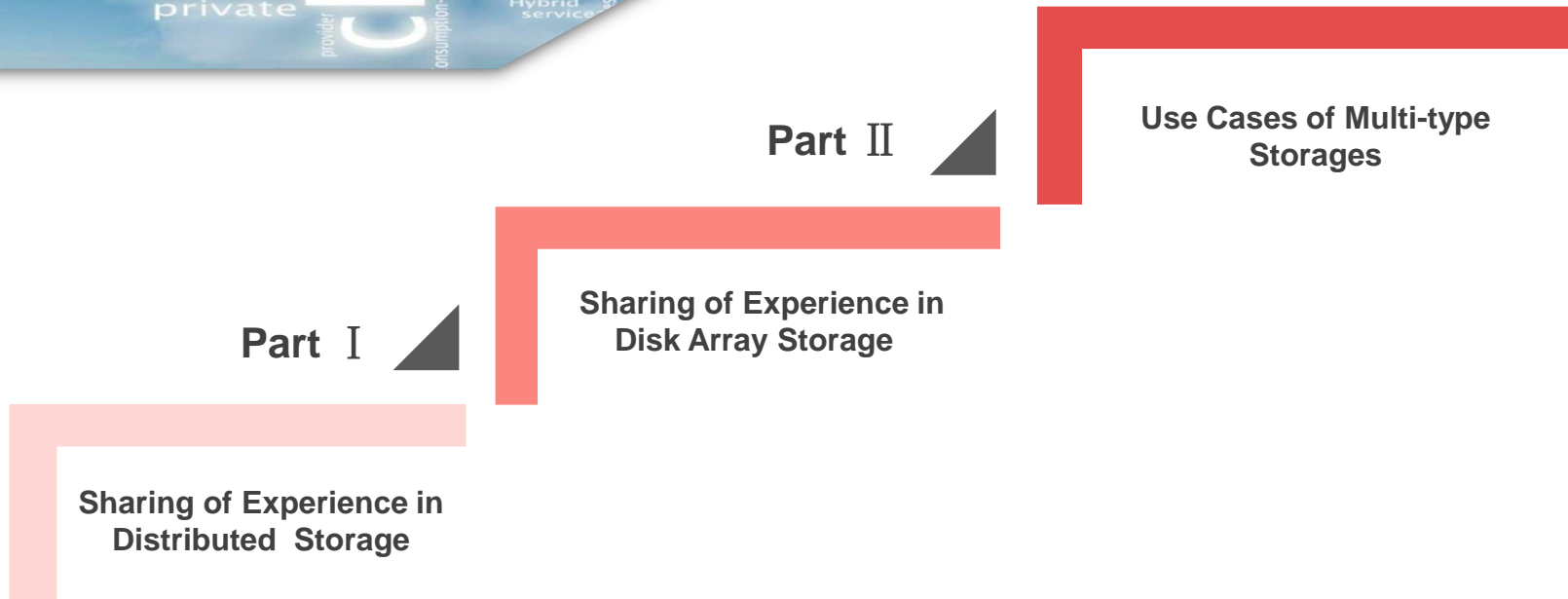
Resource type	Storage type
Critical database resource pool	Centralized storage
Vmware based resource pool	Centralized storage
KVM based resource pool	Centralized storage & distributed storage

# SRCloud Architecture and Components





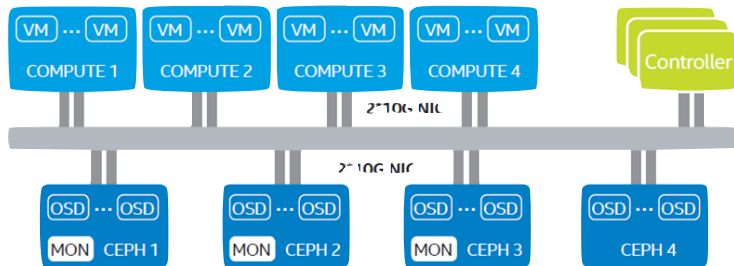
part III



# Sharing of Experience in Distributed Stor

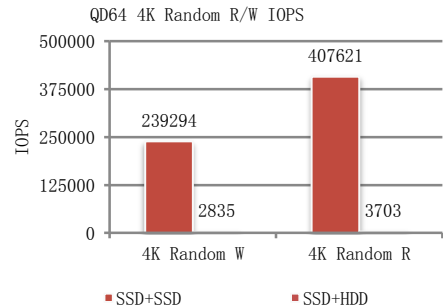
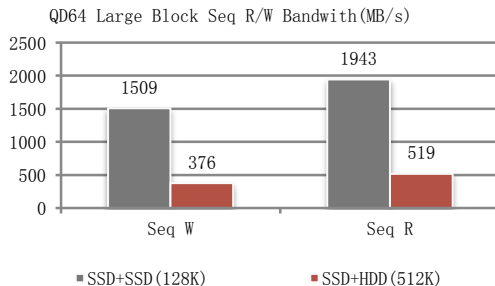
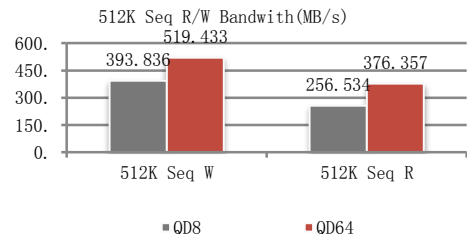
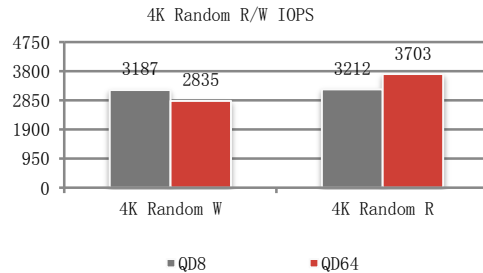
OpenStack	Ceph Monitor Nodes	Storage Nodes	OSD	Capacity
3 controller nodes 63 compute nodes	3	42	864	939T
3 controller nodes 77 compute nodes	3	39	795	879T
3 controller nodes 44 compute nodes	3	18	396	454T

## Ceph performance in the Small-scale Deployment Testing



1. SSD card + HDD disk deployment: each storage node had an Intel P3700 SSD (800G) and 8 HDDs (4T 7200 RPM), where the SSD flash memory was used for Ceph log storage and the HDDs for data storage.

2. SSD card + SSD disk deployment: each storage node had the flash memory configuration with Intel™ Optane™ P4800X SSD(375G) and Intel S3520 SSD (800G), where the Intel Optane SSD was used for log storage and the Intel S3520 SSD for data storage.



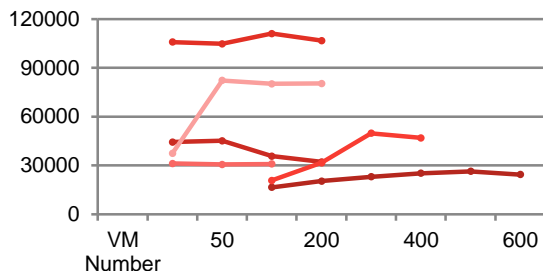
**Conclusion: in our test environment, Ceph gave full play to the performance of HDD and SSD drives.**



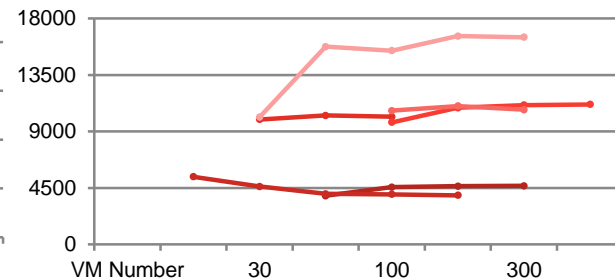
## Ceph Performance in the Mission Critical Application

Pool performance type	Journal disk type	storage nodes Number (HDD)	OSDs per node	Journal disk number	Ratio of Journal to OSD	Total OSDs
Low	HDD	27	10	N/A	N/A	270
Mid	SATA SSD	21	9	3	1:3	189
High	PCIE SSD	48	8	1	1:8	384

### 4K Random R/W IOPS



### 512K Seq R/W IOPS

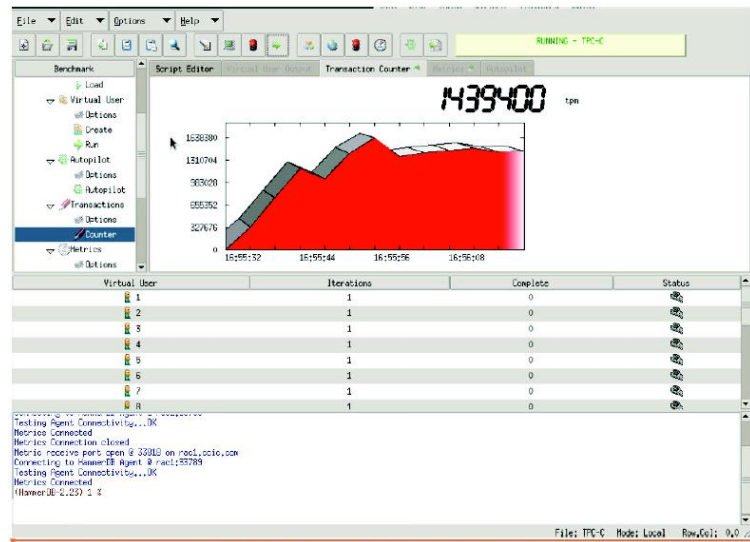
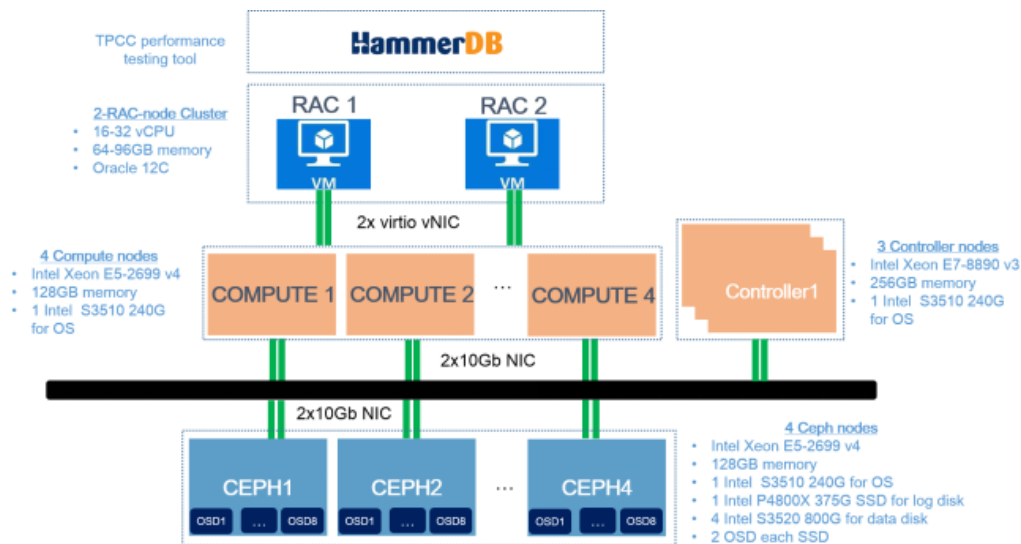


- Low Performance Pool Random W
- Middle Performance Pool Random W
- High Performance Pool Random W
- Low Performance Pool Random R
- Middle Performance Pool Random R
- High Performance Pool Random R

- Low Performance Pool Seq W
- Middle Performance Pool Seq W
- High Performance Pool Seq W
- Low Performance Pool Seq R
- Middle Performance Pool Seq R
- High Performance Pool Seq R

**Conclusion:** when the VM concurrency number of the three-type storage pools increased, the overall maximum IOPS was maintained, which could meet the operation requirements in the large-scale production environment.

## Ceph Performance in the Mission Critical Application



**Conclusion:** based on flash configuration and software optimization, Ceph can meet the performance demands of mission critical (Oracle RAC).

## Current Situation of Ceph Production Cluster on China Railway Cloud

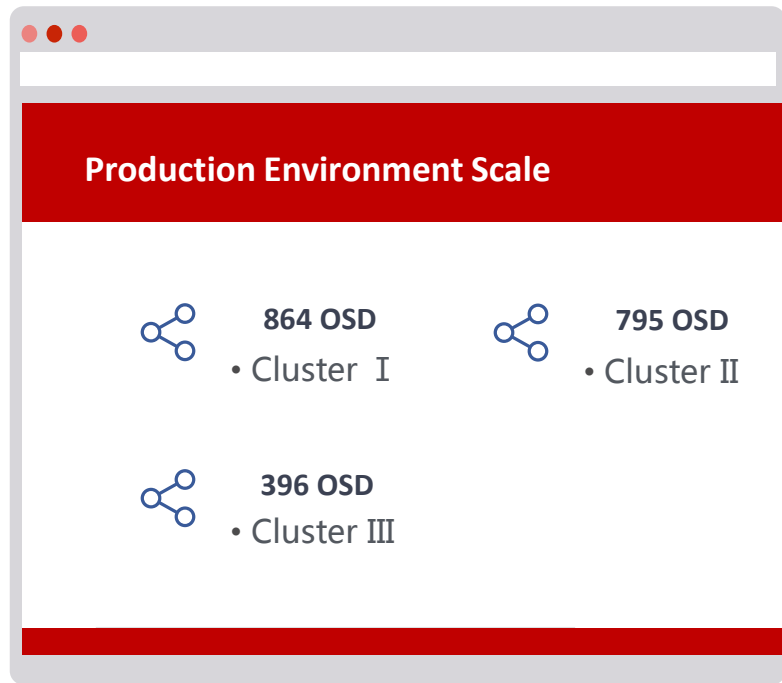
### Multi-Ceph Cluster

#### Main Operations :

- Web Server
- Application Server
- Test Database

#### Using Effect:

Ceph production cluster has been working well for over 1 year without any affect on production operation, satisfying the demands of production environment.



## Ceph 使用情况

- Pg不一致出现过20多次
- Osd所在的服务器宕机20多次
- Osd down的情况出现过4次



## Current Usage of Ceph

- Over 20 times of Pgs inconsistency
- Over 20 times of server down where Osd is running on
- 4 times of Osd down



## Three Typical Service Failures in Use

### 1.VM Performance Degradation

Two main reasons for this failure. One is that the partial OSD capacity usage exceeds the alert threshold because of non-uniform data distribution; the other is that cgroup limits osd's memory usage to ensure computing resources in the converged deployment environment, but if osd's memory is overrun, the operating system kernel will kill the osd process, leading to data migration and causing the storage cluster performance degradation.

#### **Our proposals:**

- 1.Perform the data balancing operation in the window phase, which is also added on our product management interface, based on the monitoring results after completing ceph deployment.**
- 2. In the case of converged deployment, it is suggested to adjust the cgroup limit the single osd memory to 6-8GB, rather than the recommended 2GB configuration in community.**

## Three Typical Service Failures in Use

### 2. Performance Degradation in Batch Volume Creation

- The same reason as mentioned above that cgroup has a low memory limit for osd.
- Another reason is related to the deployment architecture. If the ceph - mon service is deployed on the OpenStack control node and both share the same system disk, ceph - mon service couldn't work properly when the system disk space is insufficient.

#### **Proposal:**

**The problem should be considered before deploying ceph service. If the actual operation environment cannot satisfy the deployment of ceph - mon service separately, we should carry out the effective monitoring and warning mechanism on the hard disk.**

## Three Typical Service Failures in Use

### 3.Failure to Create Volume

Ceph pg status will often change to be inconsistent, causing the ceph cluster status to be error, so that ceph cannot provide new cloud hard disk service any more. Currently, it is unable to modify the code to fix this issue, but we do a set of effective monitoring method to conduct the timely monitor and alarm of ceph pg status to notify operation & maintenance personnel to troubleshoot as failures are happening.

**We learn from all the above problems that monitoring and alarming at the Ceph' s key indicators is so effective to master ceph besides gaining real case experience. This is a significant feature in our automated operation&management system of China Railway cloud Platform.**



## Sharing of Experience in Disk Array Storage

**In our cloud environment, we majorly use the disk array to deploy the important application systems, with the total capacity of 500T that is smaller than ceph' s storage.**

**Disk arrays from different manufacturers ( like IBM, EMC ) have different usage patterns, so that we face many troubles when managing disk array on the cloud platform. Let' s take main ones for instance.**



## 1.Failure for VM to Attach Volume

We haven't find out the root cause, but our developers' analysis suggests that it should be a bug in Linux kernel. However, rebooting the physical machine is our current solution, which may be a pragmatic and effective troubleshooting method meth in most of cases.

## 2.Disk Formatting Failure on VM and Inconsistent Disk Information

This results from lun device path on the compute node path no clearing. Our solution : First, to create a specified directory after deploying OpenStack, and then to create a soft link in the directory to the disk devices under/dev, modify the configuration of LVM. conf, and specify the LVM to scan the named directory. This is the way to avoid scanning the lun devices mounted on fc SAN.

### 3. Glance Uploading Images Slowly with Cinde As Its Back-end Storage

It was thought as a bug in glance, so we modify the logic of glance uploading images, which is that glance would create the volume based on the image's size that it is signaled before image is uploaded, so that there is only one time of disk attachment and detachment.

### 4. Failure of Volume Attachment and Detachment

This issue may result from multiple cinder-volume operating the same host on one side of the storage array. And our approach is to enable the resource pre-check. What is more, the distributed lock-management system will be introduced to settle this problem soon.

## Use Cases of Multi-type Storages

For the cloud storage usages, we not only conduct the above testings and verification, but develop many practical features by combining with both advantages and characteristics of the two storage types. Two examples:

### Cases

#### Multiple Back-end service of Storage

In a single OpenStack region, we deploy both distributed storage and disk array at the same time so that tenants can choose the more suitable storage type based on their need, and freely switch to the other one, which plays a positive role in promoting distributed storage.

#### Data Backup

We backup the data in the disk array to ceph through cinder backup, effectively reducing the backup cost of the application data.

# Practice Recommendations

- 1 Adequate testings should be conducted before the cloud services enter into production. Don't assume that the community has already done it, so there is trouble free. We have proved that lots of unexpected problems would occur in docking both distributed storage and centralized storage. So never take other project data for granted , unless your own experimental data.
- 2 To design the the size of Ceph and data disk strictly according to your actual situation.
- 3 The development of all flash memory is a revolutionary breakthrough in storage, and if possible, you can consider the full flash storage architect. Nonetheless, it still takes time to make full use of Ceph advantages and let' s work together to further improve its performance on the condition of all-flash configuration.

# Thank You

by Du Yahong

