



# 全球运维大会

2016

DevOps 2.0: 重塑运维价值



北京站

会议时间：12月16日 - 12月17日

会议地点：北京国际会议中心

主办单位：



# 大规模分布式存储开发与实践

刘源 高级架构师



# 目录

1 Sheepdog介绍

2 苏研Sheepdog应用

3 苏研Ceph应用



# 移动大云项目介绍

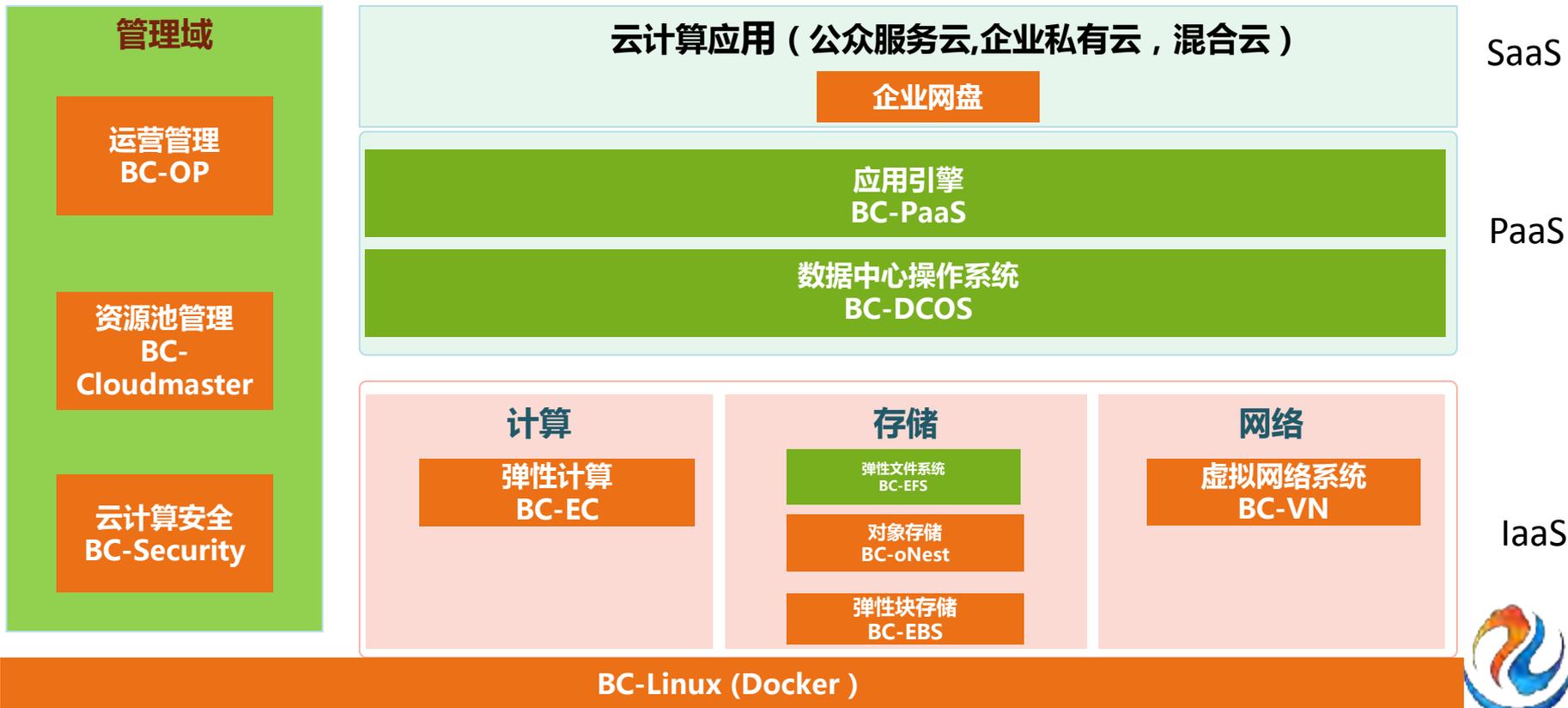
计算



存储



# 移动云计算产品



# Sheepdog介绍

- 用户态分布式对象存储系统

- 管理磁盘和节点

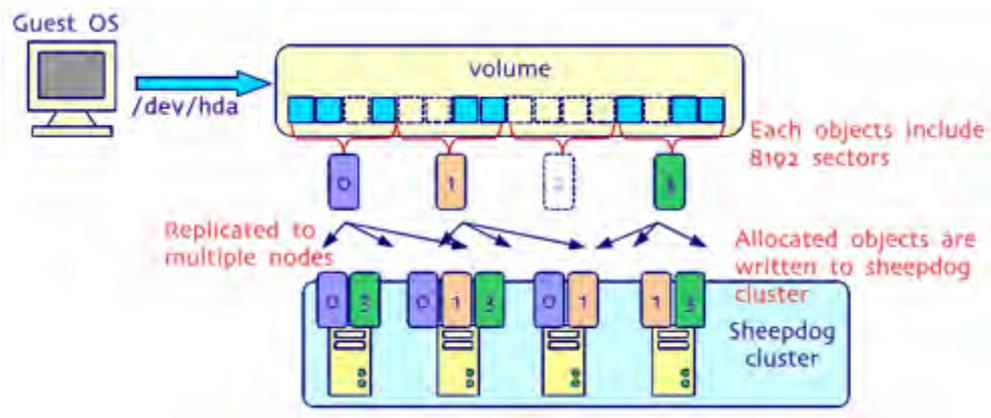
- 聚合容量和能力(IOPS + 吞吐量)
    - 隐藏硬件故障
    - 动态扩容或缩容

- 保护数据

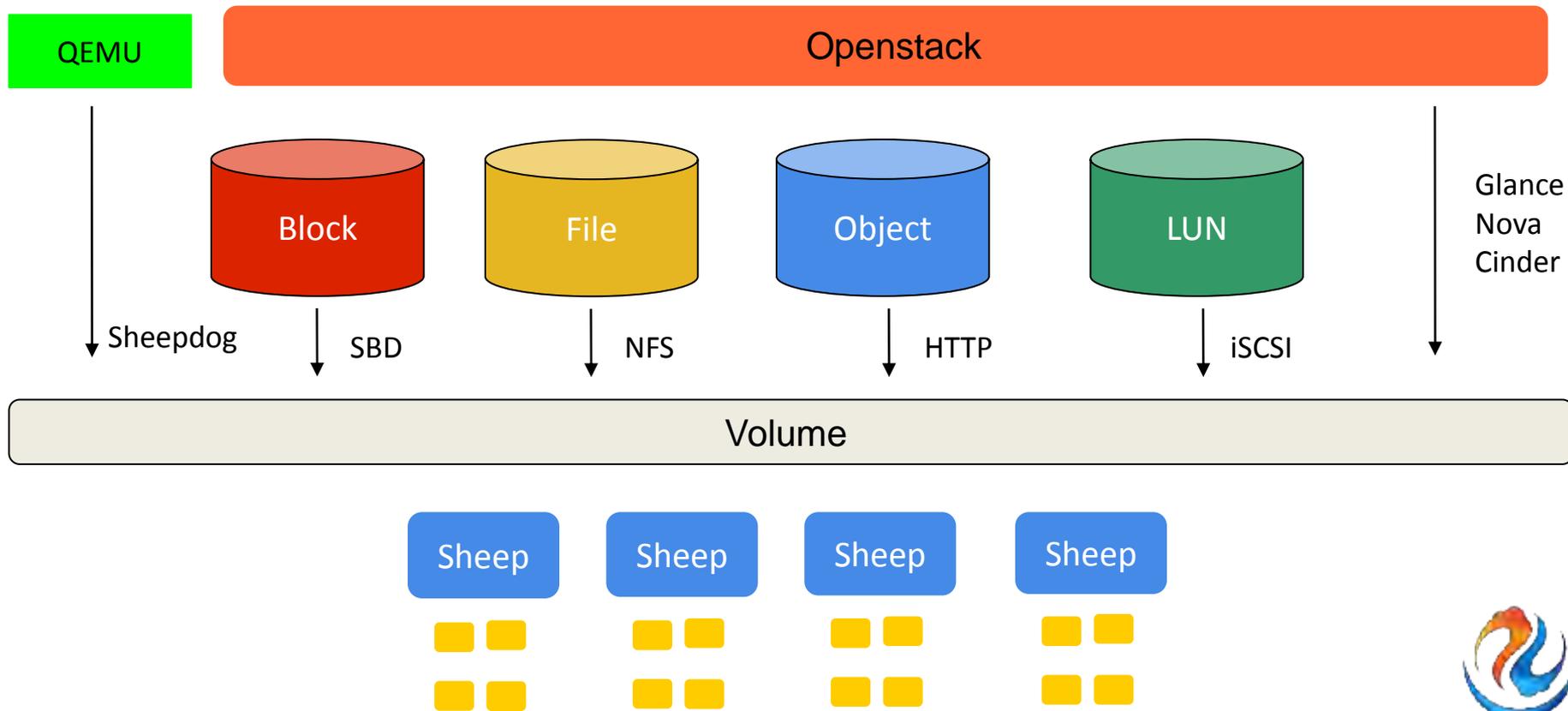
- 提供冗余机制(复制和纠删码)
    - 提供数据自动迁移和修复机制

- 多样化接口 (单个集群)

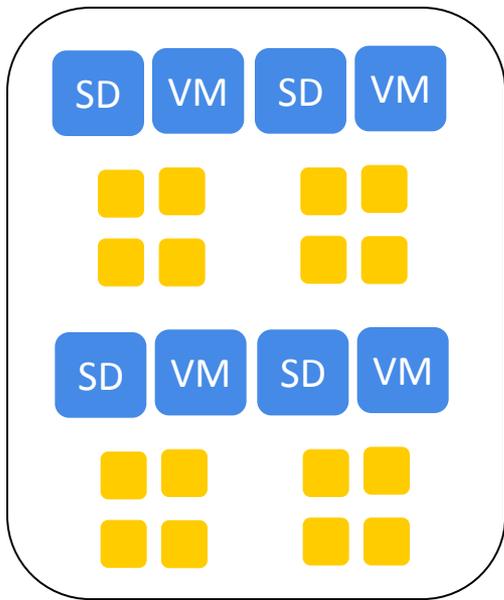
- 虚拟机虚拟卷
    - 对象存储(Swift和S3协议)
    - Openstack (Cinder, Glance, Nova) 原生支持
    - 通过NFS协议提供文件接口
    - Linux内核级块设备



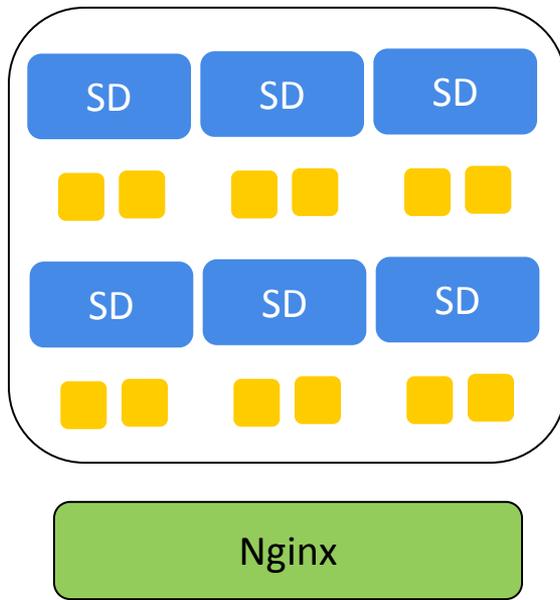
# 软件全栈视图



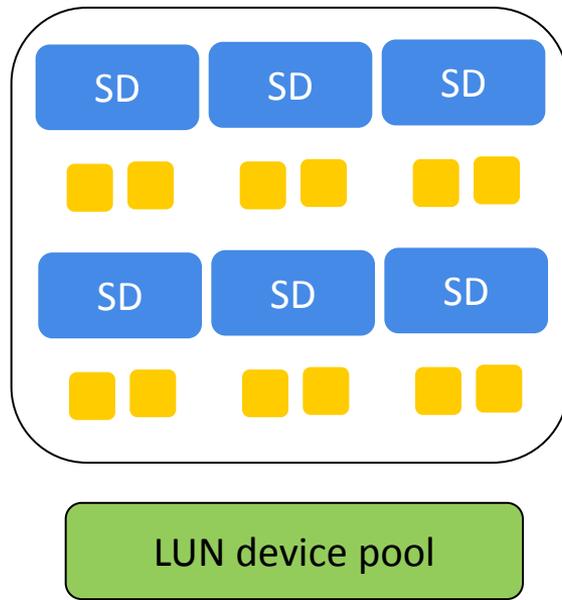
# 应用方案



计算存储融合



对象存储



计算存储分离



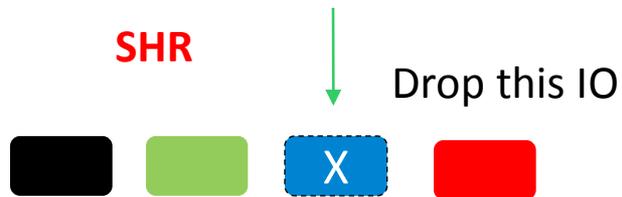
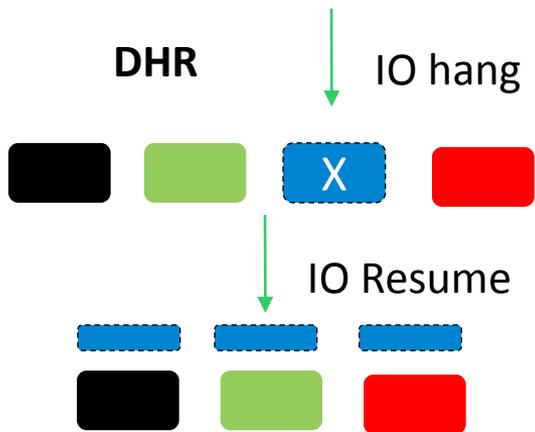
# Sheepdog上线情况

- 500+节点
  - 单集群10~40个节点，单节点12块机械硬盘，计算存储分离
- 服务
  - STGT + sheepdog: 通过iSCSI LUN提供SAN服务
  - QEMU 直连：提供虚拟机块服务



# 线上问题

- 大多数都是网络故障（网卡down，交换机问题，网络抖动等），少量磁盘故障
  - 建议：Zookeeper最好定时重启
- 集群恢复时服务质量下降严重
  - 解决方案：开发静态哈希特性，避免单节点故障触发全集群恢复



- 80%的故障都是单节点故障
- 数据有冗余，短期内不用恢复
- 单节点恢复，而不是集群恢复
- 避免集群恢复



# 线上问题

- 缺乏有效监控，运维效率低
- 系统升级内核以及glibc Bug导致数据异常，应用Crash
  - 建议：不要贸然升级操作系统

## Bug 1293976 - CVE-2015-5229 glibc: calloc() returns non-zero'ed memory [rhel-7.3.0]

Status: CLOSED ERRATA

Aliases: None

Product: Red Hat Enterprise Linux 7  
Component: glibc (Show other bugs)  
Version: 7.2  
Hardware: All Linux

Reported: 2015-12-23 16:05 EST by Jeff Layton

Modified: 2016-11-03 04:28 EDT (History)

CC List: 17 users (show)

See Also:

Fixed In Version: glibc-2.17-107.el7

Doc Type: Bug Fix



# 线上问题

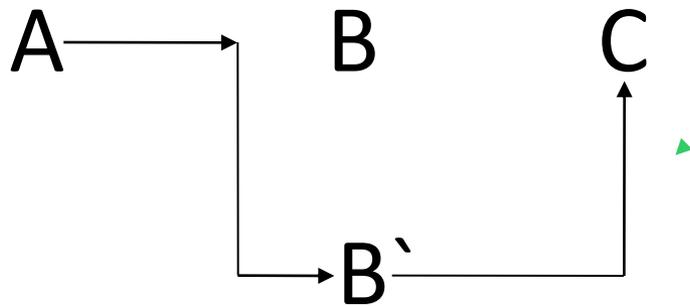
- 不能热升级

- 解决方案：热补丁

- 原理：0 二进制代码提取 1 动态加载(类似.so) 2 函数替换 (类似内核Ftrace)
- 通常对小问题修复很有效，不能修改数据结构和全局变量。

A -----> B -----> C **After Patching**

```
diff --git a/sheep/cluster/zookeeper.c b/sheep/cluster/zookeeper.c
index 487455d..c3f9f8a 100644
--- a/sheep/cluster/zookeeper.c
+++ b/sheep/cluster/zookeeper.c
@@ -1209,6 +1209,7 @@ static const int zk_max_event_handlers = ARRAY_SIZE(zk_event_handlers);
 static inline void handle_session_expire(void)
 {
     /* clean memory states */
     close(efd);
     zk_tree_destroy();
     INIT_RB_ROOT(&zk_node_root);
 }
```



B' is loaded by Linux's dynamic loader on the fly



# 转型Ceph

- 主要原因
  - Ceph越来越稳定
  - 虽然复杂，但社区非常活跃，问题快速响应
  - Openstack支持度最好
  - 厂商支持（开发运维监控）

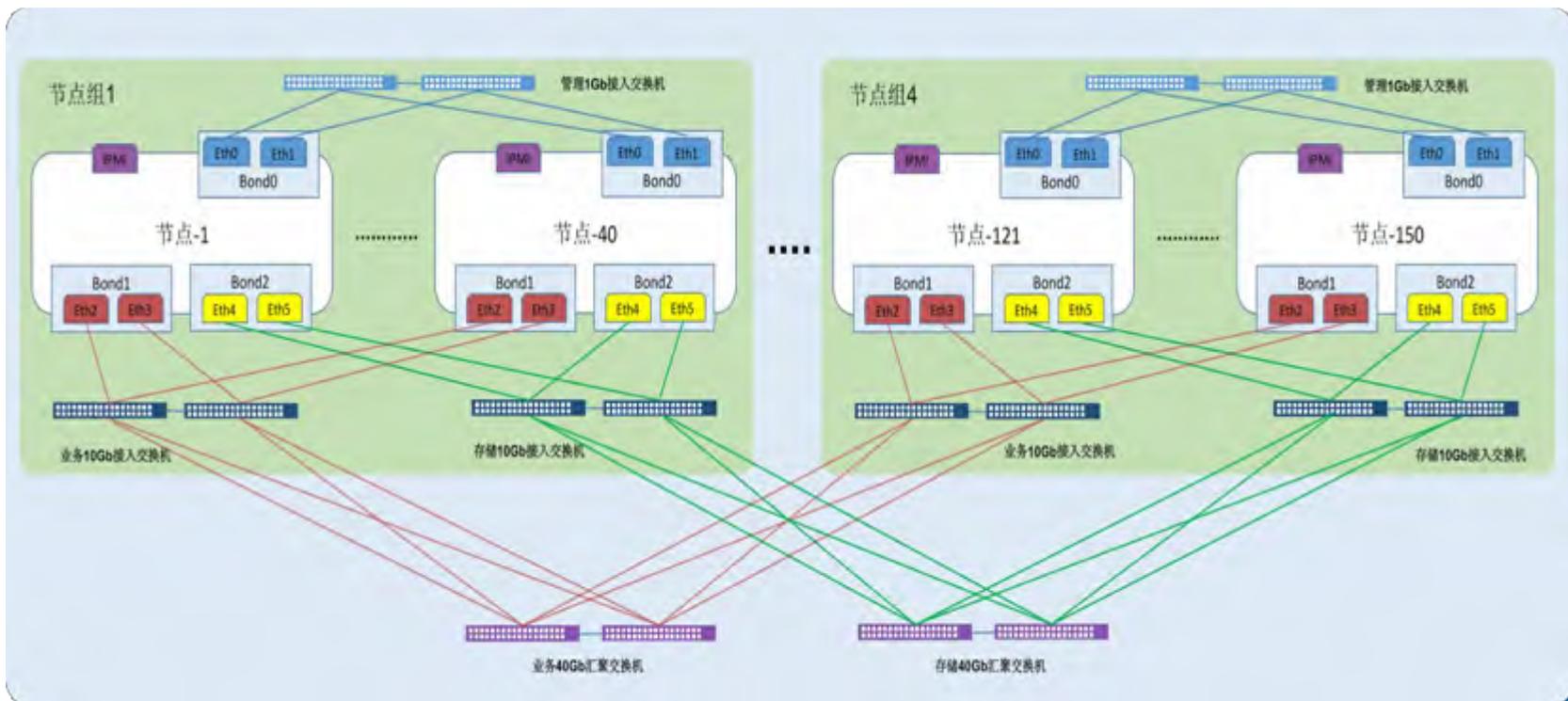


# 苏研Ceph上线情况

- 基于J版
- 对象存储
  - 上线500+节点
  - 单集群100+节点，单节点24个硬盘，可能是全球最大的集群(2k+ OSD)
  - 可能是全球首个异地容灾生产案列
- 块存储
  - 虚拟机块设备 (librbd)和iscsi LUN服务
  - NAS服务 (librbd + nfsd in VM)
  - 上线100+节点
  - 小集群部署

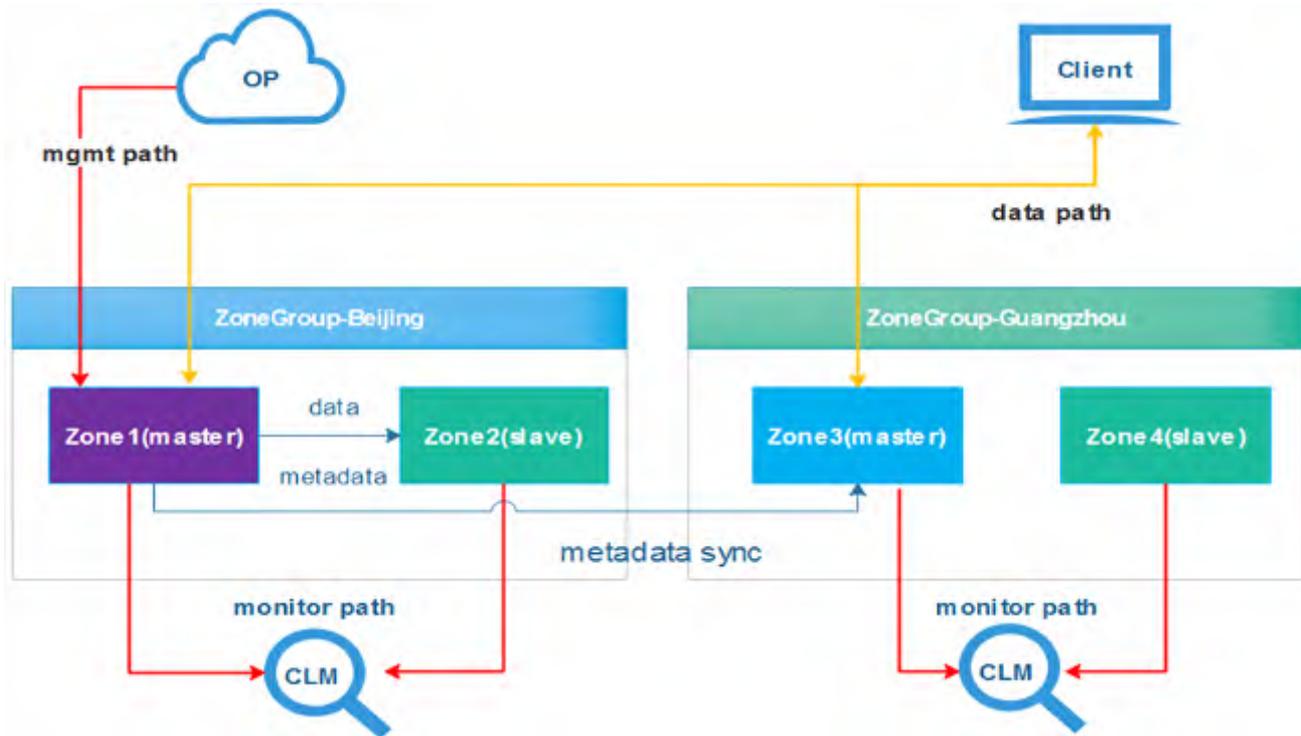


# Ceph部署架构图



# Ceph多数据中心应用

两个数据中心，异地备份



# 线上问题

- QEMU Centos 7系列虚拟机 virtio\_blk性能很差
  - virtio\_scsi性能正常
  - 原因：centos7.1内核的bdev文件系统未实现readpages接口
  - 建议线上部署多用virtio\_scsi而不是virtio\_blk，这也是红帽推荐的



## 苏研存储技术号



# DevOpsDays 即将首次登陆中国



DevOps 之父 Patrick Debois 与您相约  
DevOpsDays 北京站 2017年3月18日



门票早鸟价仅限前100名，请从速哟

<http://2017-beijing.devopsdayschina.org/>



GOPS2016  
Beijing



想第一时间看到  
高效运维社区公众号  
的好文章吗？

请打开高效运维社区公众号，点击右上角小人，如右侧所示设置就好





# Thanks

高效运维社区  
开放运维联盟

荣誉出品

