



# Bring SCSI support into QEMU block layer

白耀伟

baiyaowei@cmss.chinamobile.com

www.10086.cn

- Backgroud
- Solution
- Howto
- Status
- Future work

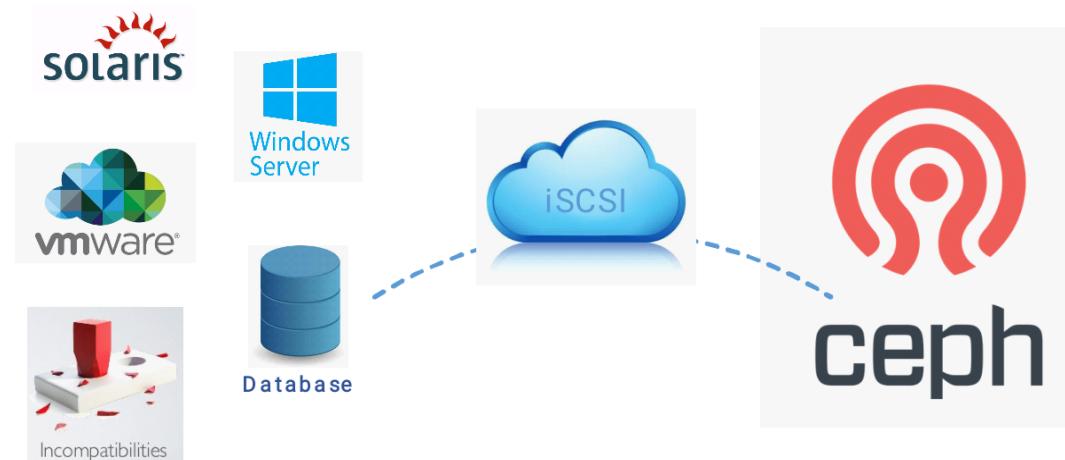


# Background

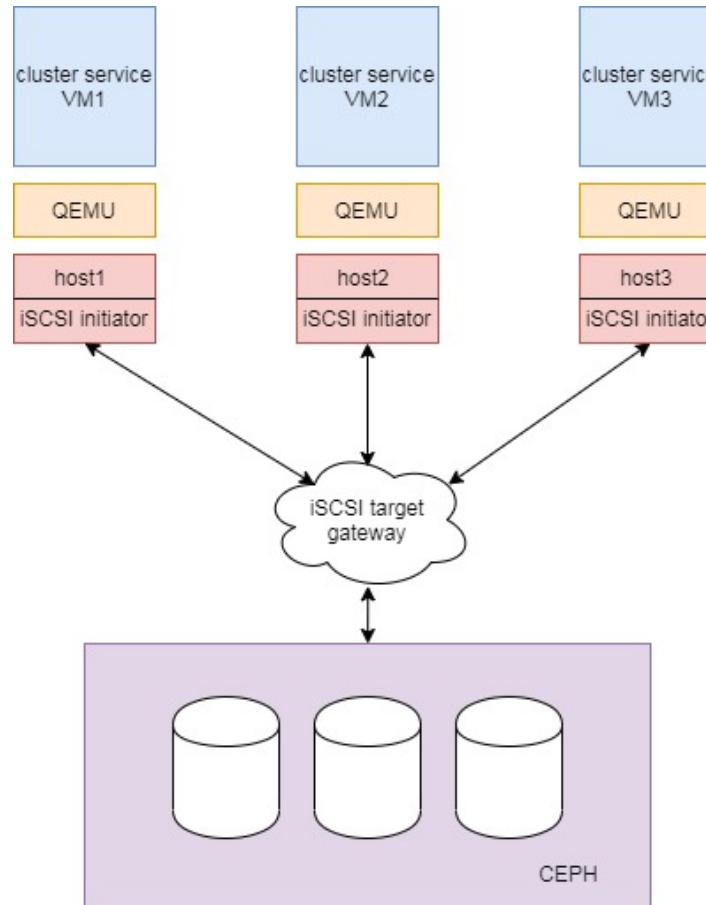
[www.10086.cn](http://www.10086.cn)

# Background

- Clustering services, like OCFS, MSCS, share disks with concurrency control mechanisms
- Concurrency control mechanisms are usually implemented via block-level protocols, like iSCSI
- Shared disks are supplied by distributed block storage, like Ceph



# Background



# Background

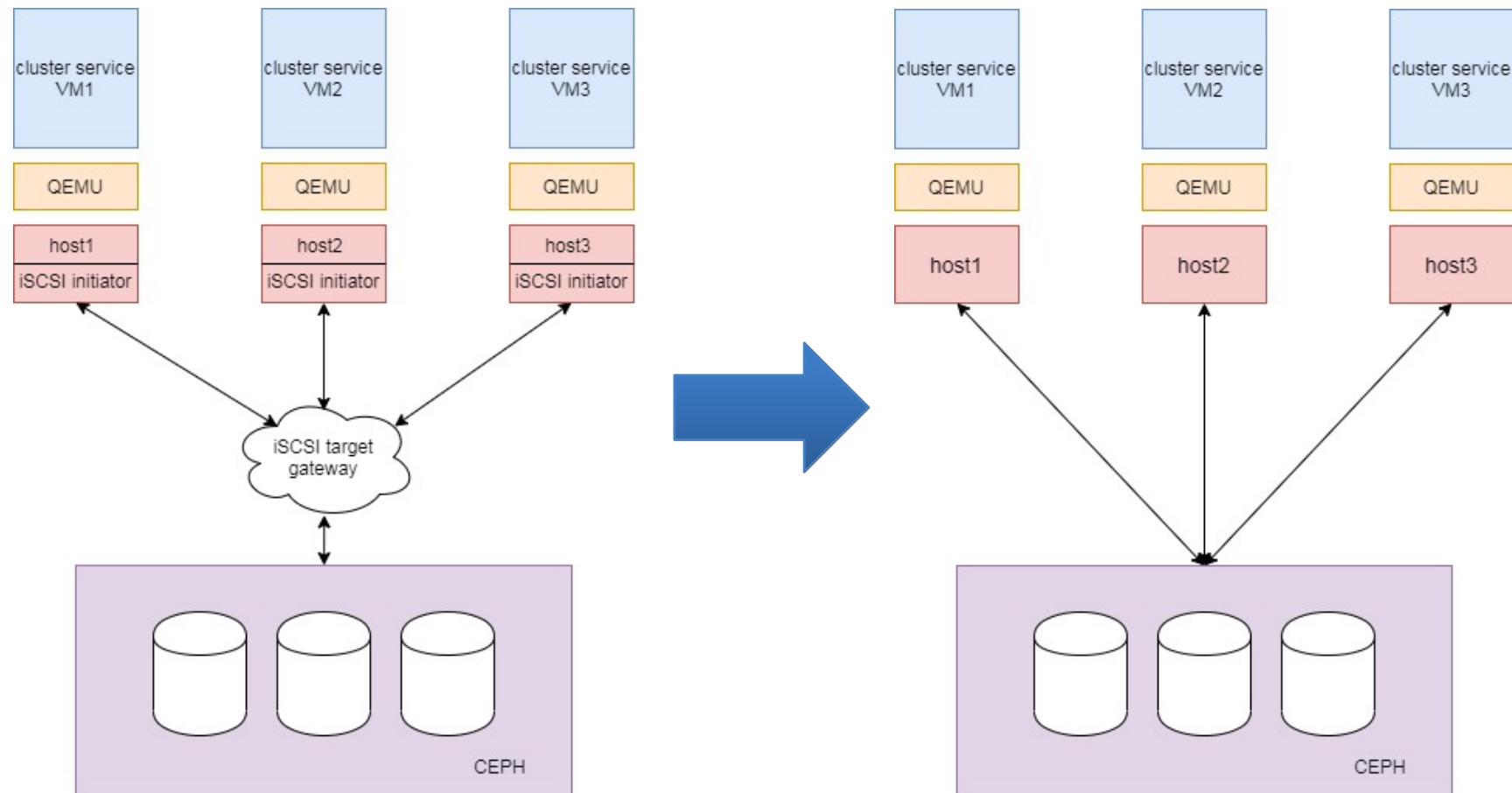
- Problems:
  - Long IO path
  - Immature components
  - Hard to maintain



# Solution

[www.10086.cn](http://www.10086.cn)

# Solution



- Work to do
  - SCSI support in Ceph(done)
    - COMPARE AND WRITE(upstream)
    - PERSISTENT RESERVATION(private)
  - SCSI support in QEMU(miss)
    - COMPARE AND WRITE
    - PERSISTENT RESERVATION



# Howto

[www.10086.cn](http://www.10086.cn)

- SCSI support in QEMU
  - SCSI device emulation
  - Block layer interface
  - Block IO path interface
  - Block driver interface

- SCSI support in QEMU
  - SCSI device emulation
    - hw/scsi/scsi-disk.c
    - COMPARE AND WRITE, PERSISTENT RESERVATION emulation
  - Block layer interface
  - Block IO path interface
  - Block driver interface

- SCSI support in QEMU
  - SCSI device emulation
  - Block layer interface
    - block/block-backend.c
    - reuse blk\_aio\_pwritev for COMPARE AND WRITE
    - new blk\_persistent\_reserve\_{in,out,check}
  - Block IO path interface
  - Block driver interface

- SCSI support in QEMU
  - SCSI device emulation
  - Block layer interface
  - Block IO path interface
    - block/io.c
    - reuse bdrv\_driver\_pwritev for COMPARE AND WRITE
    - new bdrv\_persistent\_reserve\_{in,out,check}
  - Block driver interface

- SCSI support in QEMU
  - SCSI device emulation
  - Block layer interface
  - Block IO path interface
  - Block driver interface
    - block/rbd.c
    - new bdrv\_aio\_COMPARE\_and\_WRITE and bdrv\_co\_persistent\_reserve\_{in,out,check}



# Status

[www.10086.cn](http://www.10086.cn)

- Online in our IaaS cloud service
- Patchset(COMPARE AND WRITE part)
  - :
  - <https://patchwork.ozlabs.org/project/qemu-devel/list/?series=141595>



# Future work

[www.10086.cn](http://www.10086.cn)

- Upstream
  - CEPH
    - PERSISTENT RESERVATION
  - QEMU
    - WRITE SAME
    - PERSISTENT RESERVATION



# Q&A

[www.10086.cn](http://www.10086.cn)