

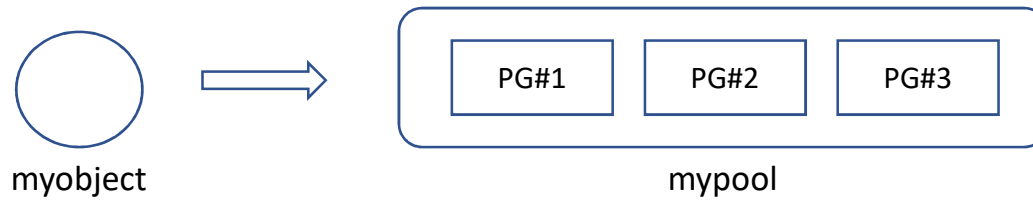
Ceph & RocksDB

변일수(Cloud Storage팀)

© LINE

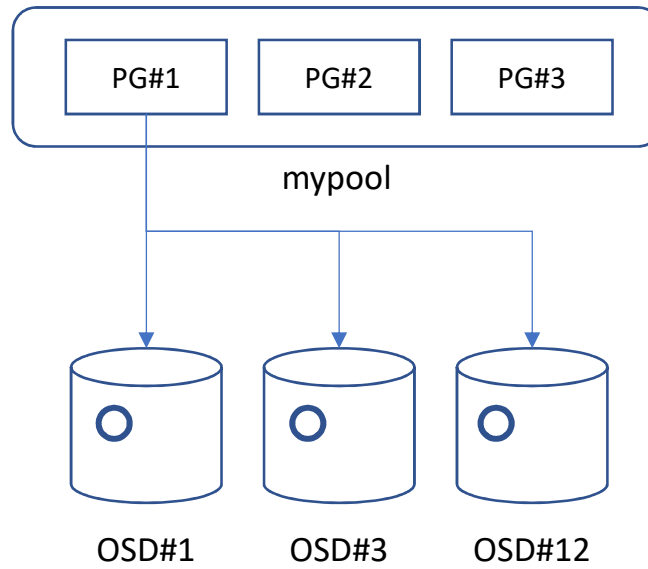
Ceph Basics

Placement Group

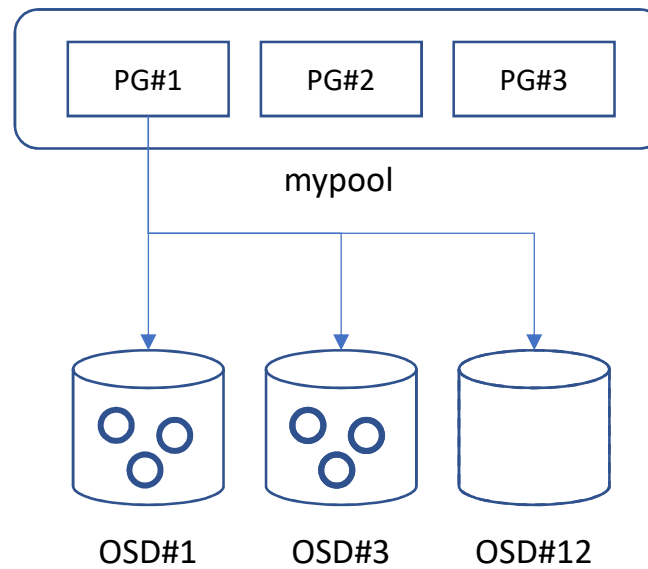


$$\text{hash(myobject)} = 4 \% 3(\# \text{ of PGs}) = 1 \leftarrow \text{Target PG}$$

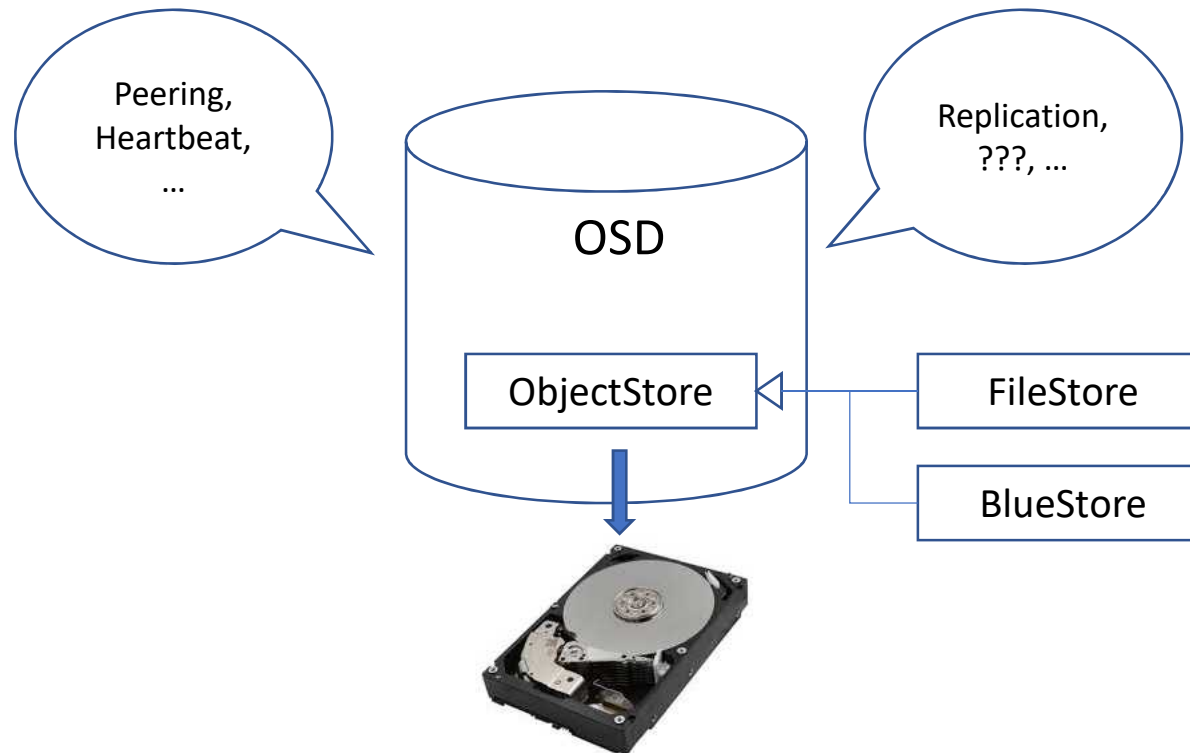
CRUSH



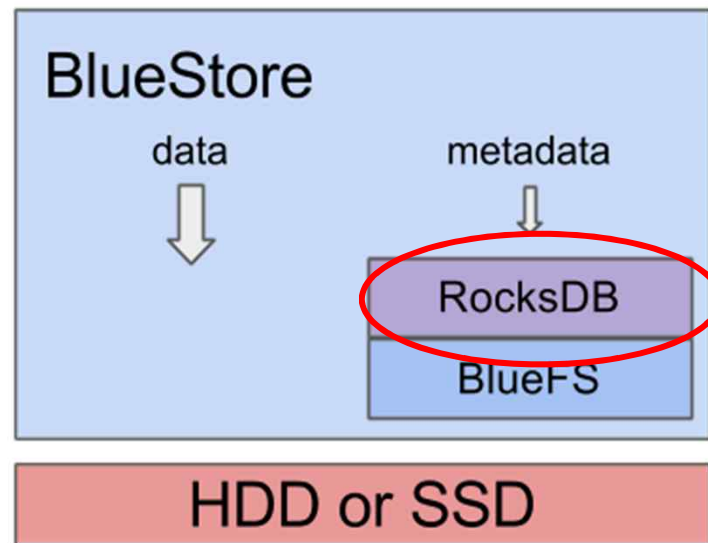
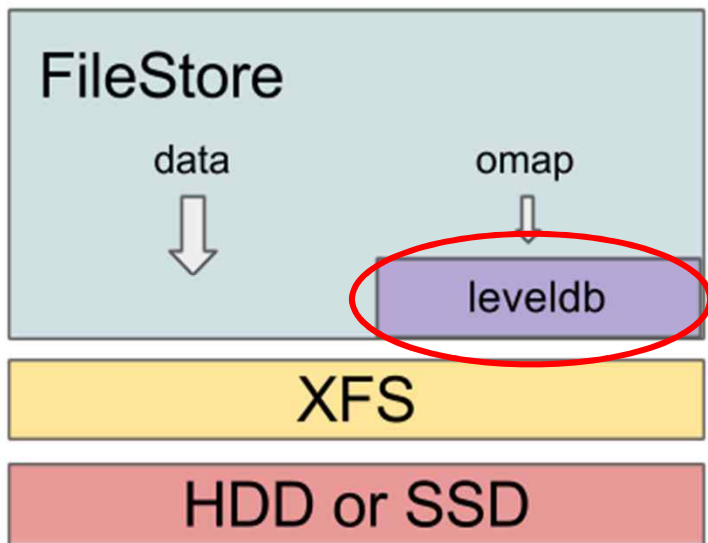
Recovery



OSD



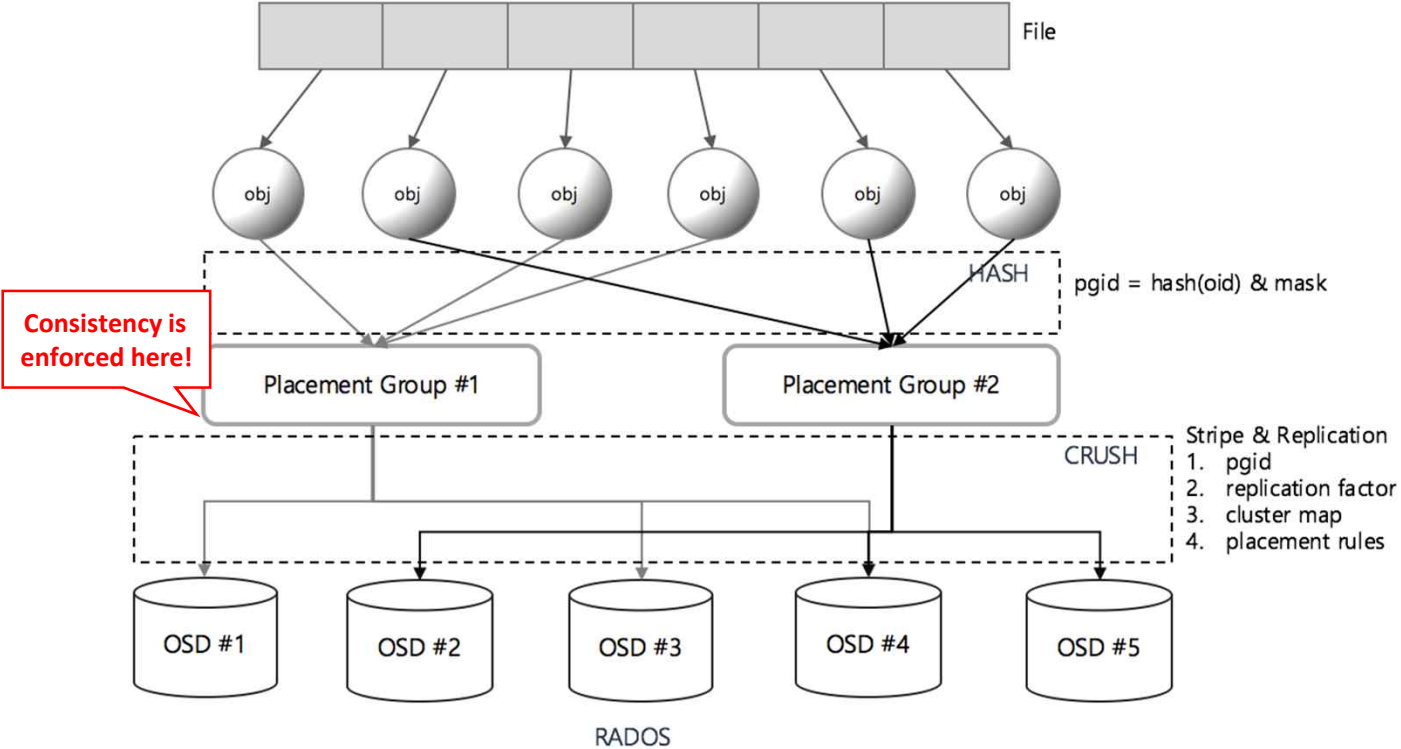
ObjectStore



<https://ceph.com/community/new-luminous-bluestore/>

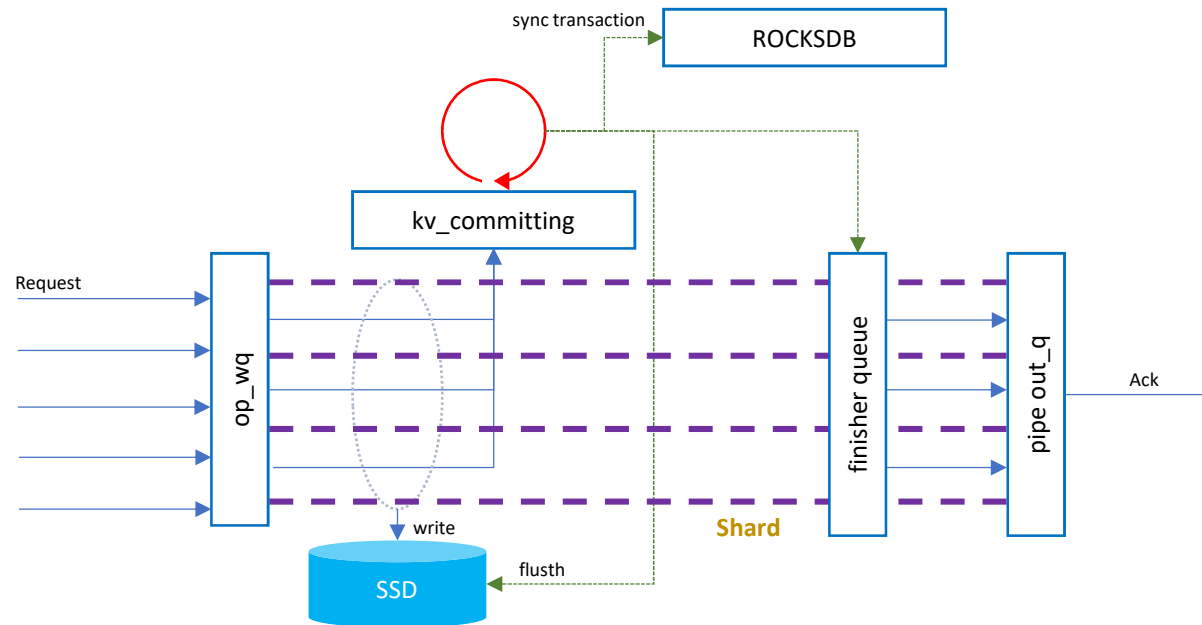
OSD Transaction

CRUSH



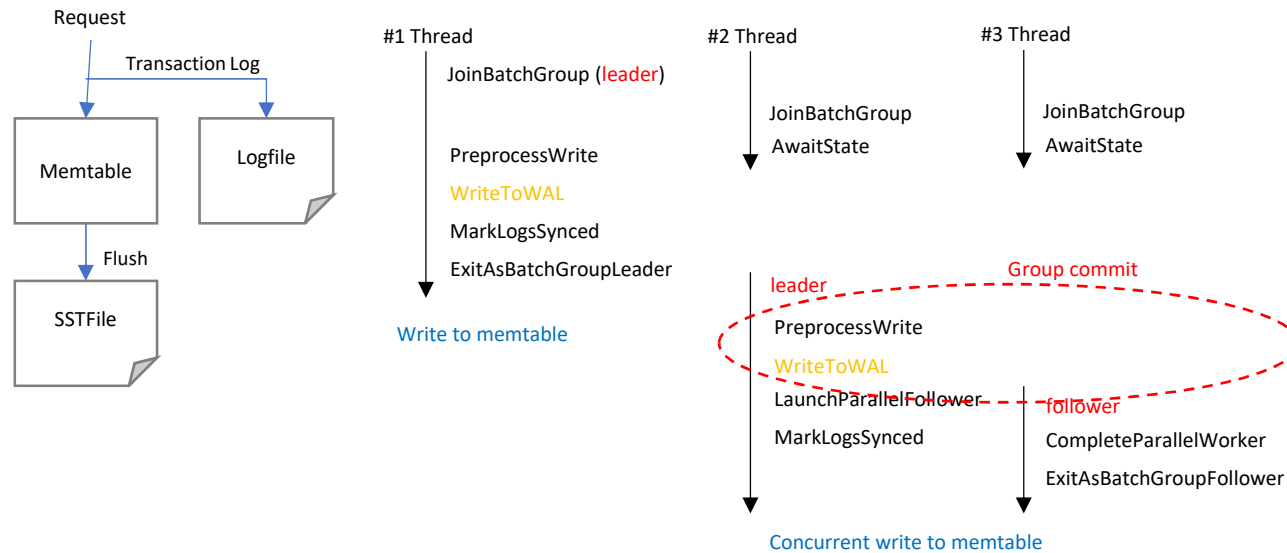
BlueStore Transaction

- To maintain ordering within each PG, ordering within each shard should be guaranteed.

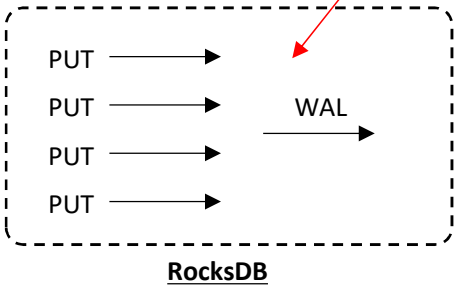
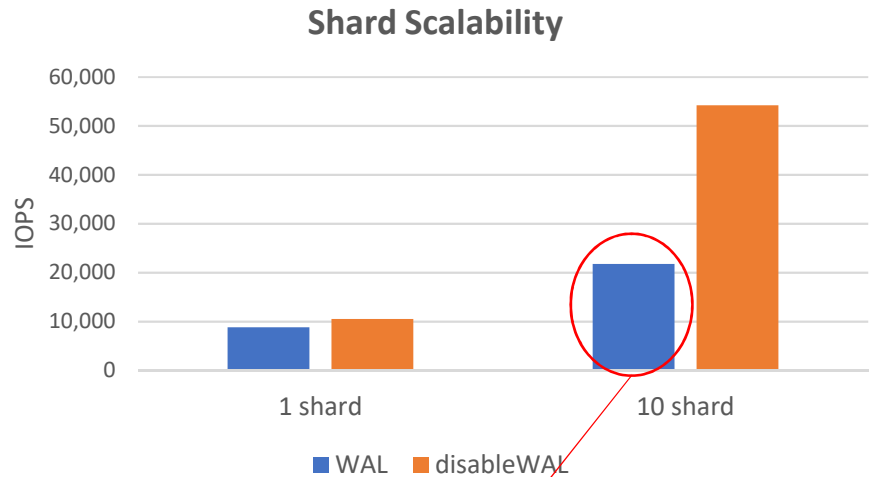


RocksDB Group Commit

- Metadata is stored in RocksDB.
- After storing metadata atomically, data is available to users.



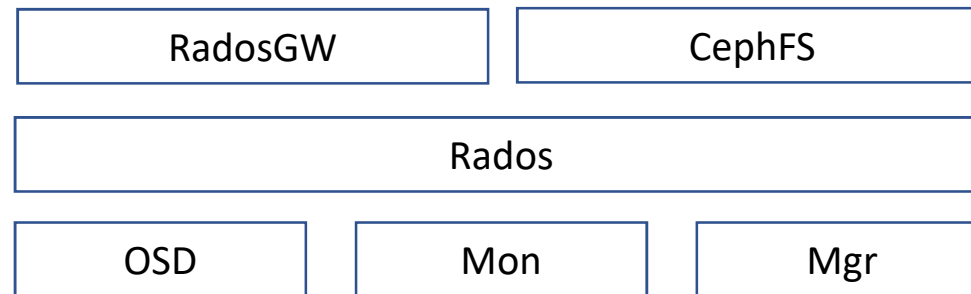
Thread Scalability



RadosGW

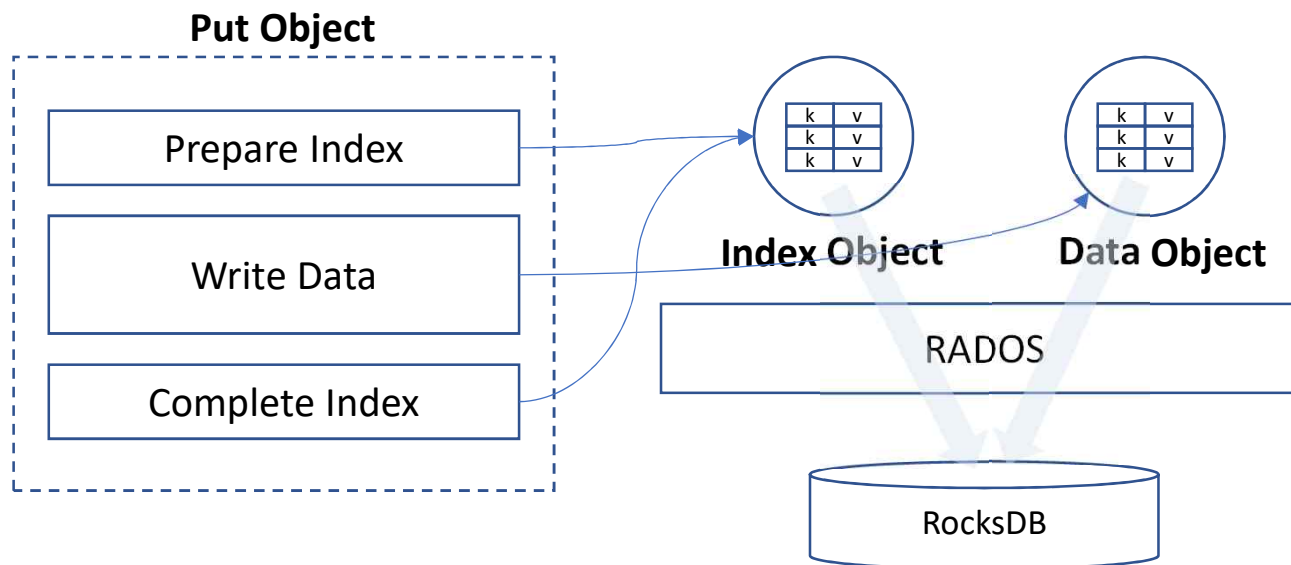
RadosGW

- RadosGW is an application of RADOS



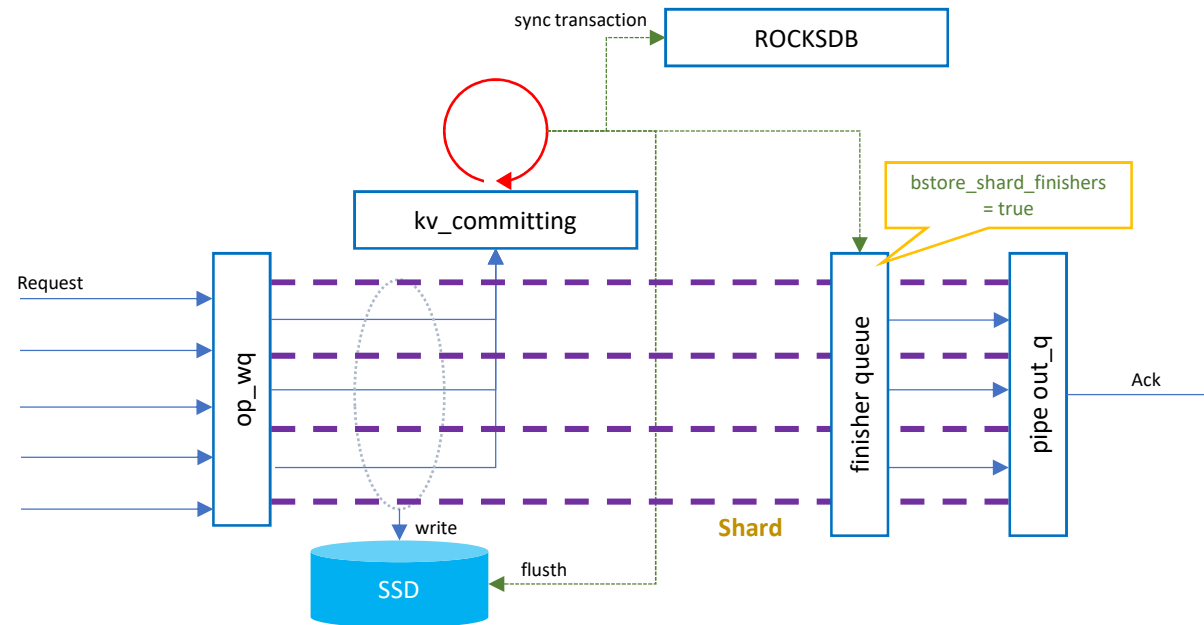
RadosGW Transaction

- All atomic operations depend on RocksDB



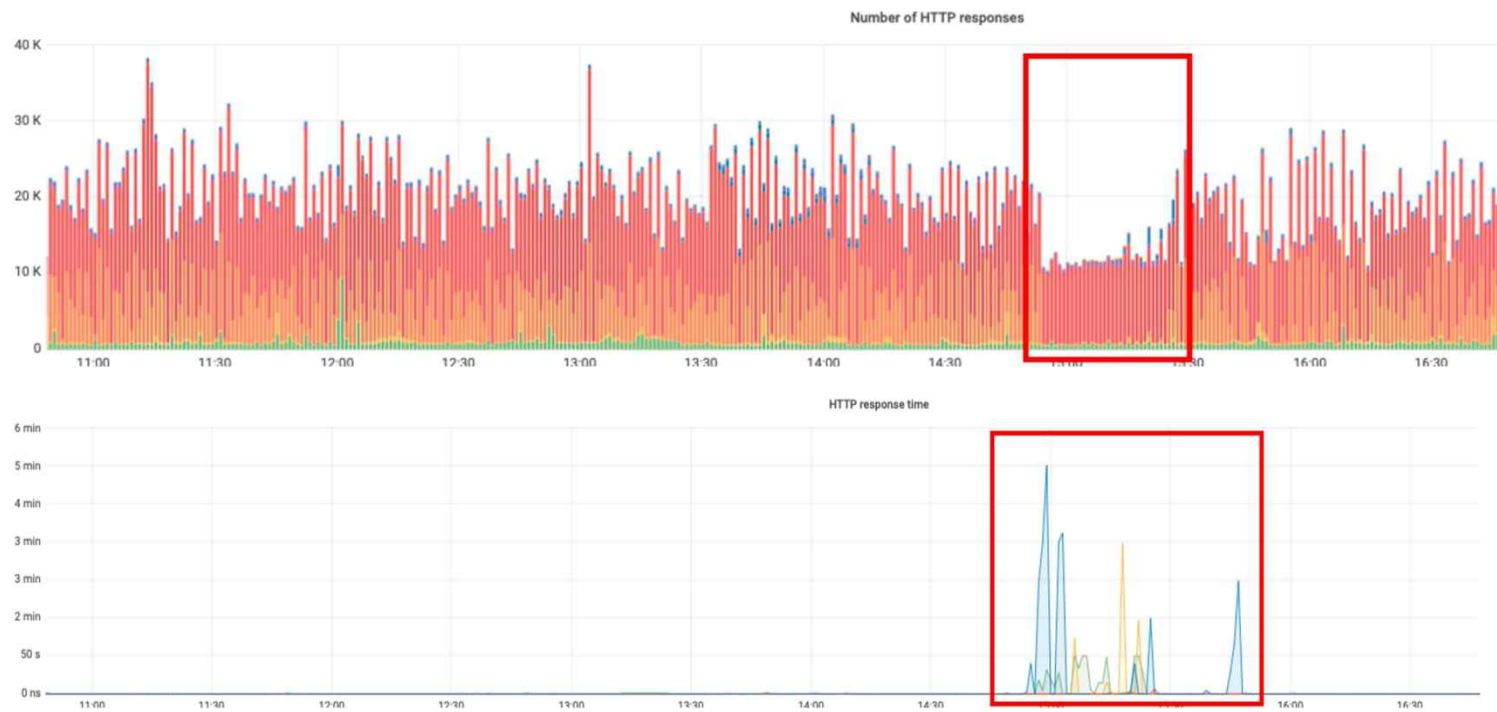
BlueStore Transaction

- To maintain ordering within each PG, ordering within each shard should be guaranteed.

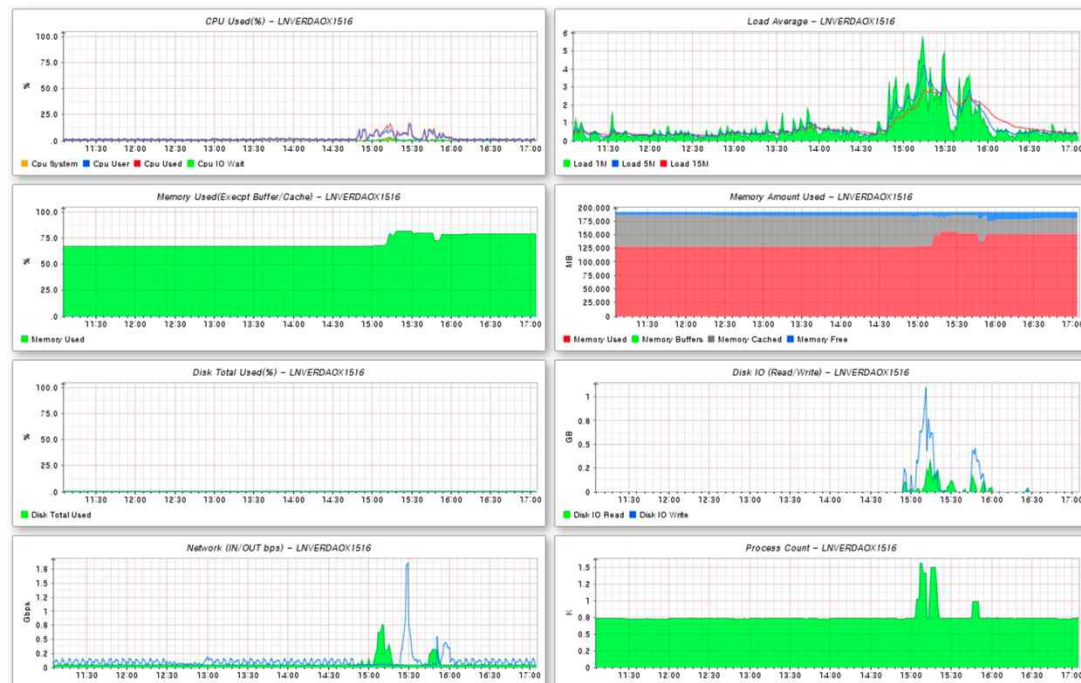


Performance Issue

Tail Latency

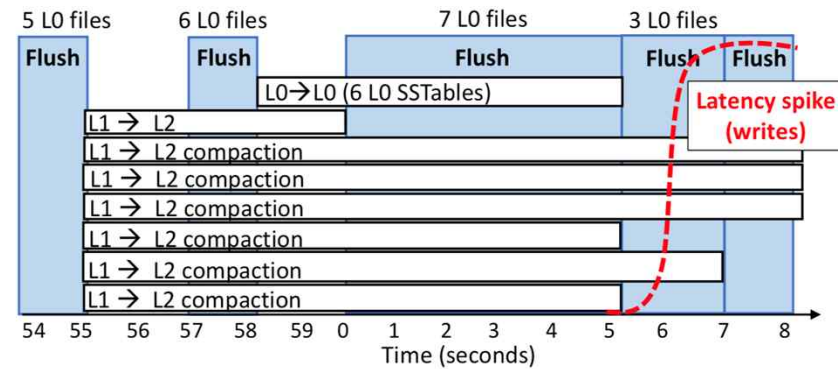
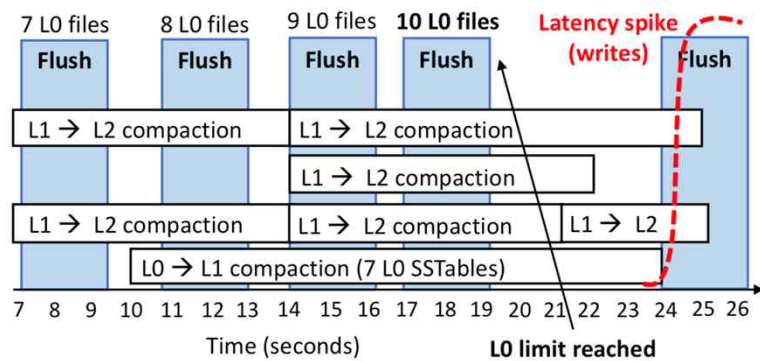


Performance Metrics



RocksDB Compaction Overhead

- "[SILK: Preventing Latency Spikes in Log-Structured Merge Key-Value Stores](#)" (ATC'19)



Conclusions

- Ceph highly depends on RocksDB
 - Strong consistency of Ceph is implemented using RocksDB transactions
- The performance of ceph also depends on RocksDB
 - Especially for Small IO
- But RocksDB has some performance Issues
 - Flushing WAL
 - Compaction
- ilsoobyun@linecorp.com

THANK YOU