hosted by **Alibaba** Group 阿里巴巴集团  APACHE HBASE

# HBaseConAsia2018

Gehua New Century Hotel Beijing,  China

August 17,2018

# HBase On Persistent Memory

Anoop Sam John,  Ramkrishna S Vasudevan

(intel)

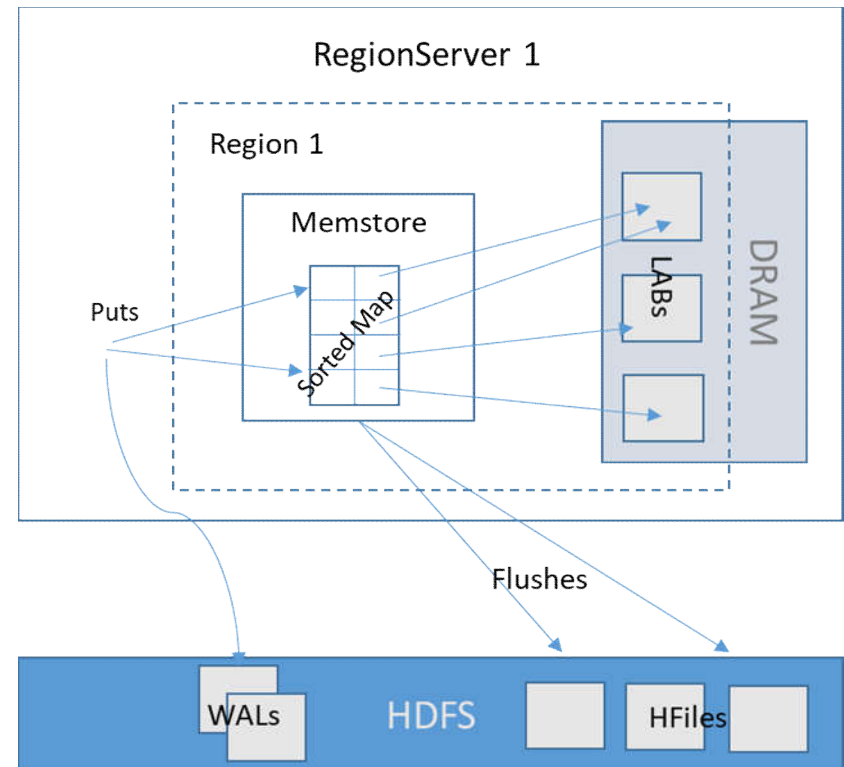# Content

hosted by **Alibaba** Group 阿里巴巴集团   **HBASE**

# Apache HBase Present Model

❖ Accumulate data in Memory

o Sorted Map

o Cell data bytes in Local Allocation Buffers (LABs) in DRAM

o LABs with size 2 MB in RAM.

❖ Also write to Write Ahead Log (WAL)

o Cell data in Volatile RAM.

o To recover from server crash

o HDFS interaction adding more latency

　hsync vs hflush - HBASE-19024

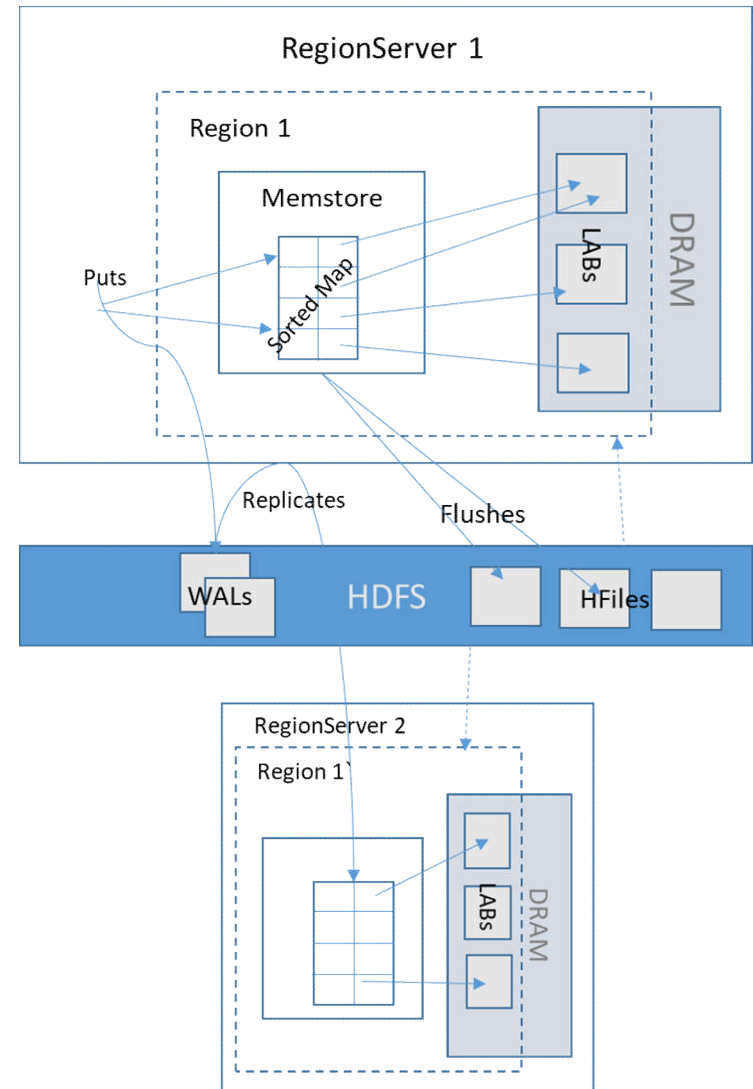❖ Flushes as files to HDFS on reaching memstores size

o 128 MB default flush size

❖ Replay WAL on server crash

o Data unavailable till replay completes

o Large Mean Time To Recover (MTTR)

o Takes several minutes on large cluster *(Complaint from many*

　*users like Alibaba – HBaseConAsia , Huawei)*

# Apache HBase Present Model
Region Replica for better availability (?)

❖ Read only replica regions on other RSs

- o Refer to same HFiles in HDFS

- o Memstore data also can be replicated

  - o Using WAL read replication path

  - o Eventual consistency

- o Can read from replica regions when primary is down

  - o No strong consistency guarantees.

  - o Only when Scan/Get says TIMELINE Consistency, replica read happens (not by default)

- o Better availability but only for selected use cases !!! – No strong Consistency

# Persistent Memory Technology

## Persistent Memory

- o Get back data even after power cycles
- o Accessed using memory APIs
- o Processor load and store instructions

## Intel Apache Pass (AEP)

- o Persistent Memory (pmem)
- o Big, affordable and persistent
- o Accessed like volatile memory, using processor load and store instructions
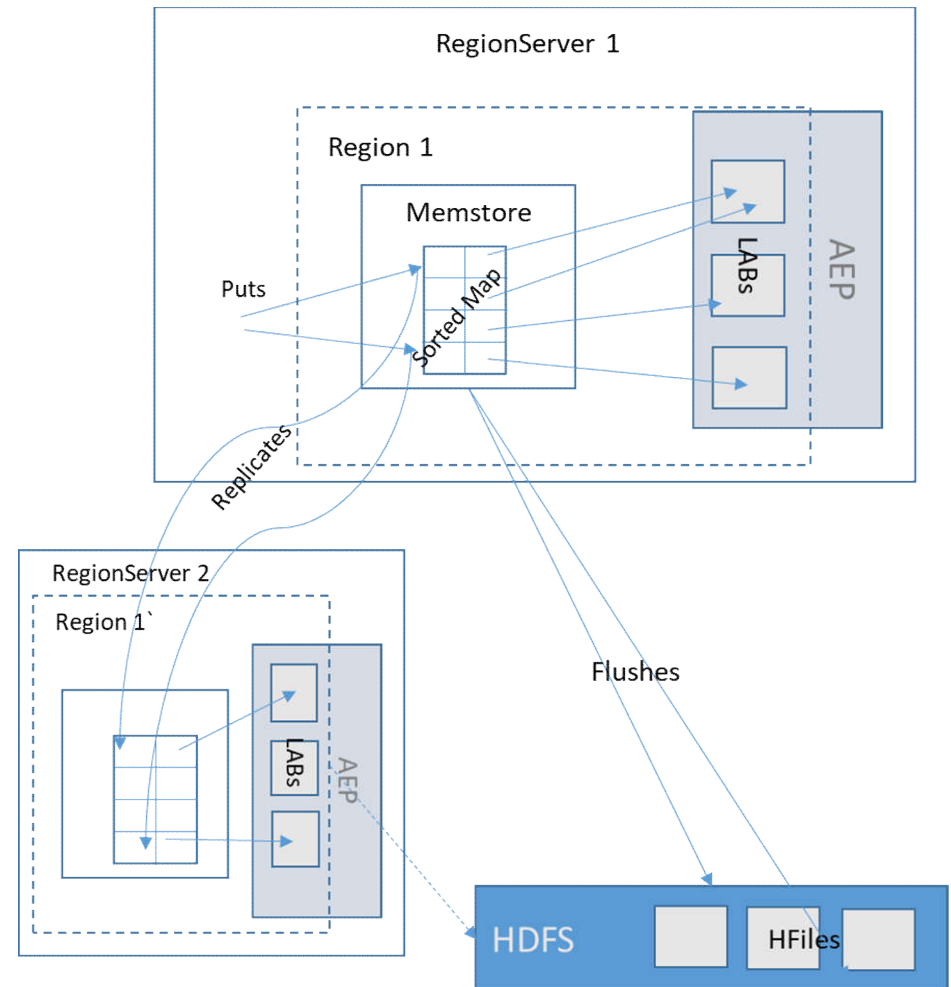
## 3DXPoint

- o New NVM technology
- o Stackable cross-gridded data access array

## Library

- o NVML (Now called PMDK)
- o Java wrappers around it
- o    Apache Mnemonic – (used for PoC)
- o    *https://github.com/pmem/pcj*

# Apache HBase On Persistent Memory

❖ Accumulate data in Memory

o Sorted Map

o Cell data bytes in Local Allocation Buffers with size 2 MB

❖ Region Replica in other servers

o Replica regions feature already in place.

o Synchronous replication to replica regions

❖ No need to write to Write Ahead Log (WAL)

o Cell data in non volatile AEP.

❖ Server crash

o Fast switch to replica regions

o Consistent data in replica regions

o Full cluster down – Data in non volatile area. Fast way to recreate in memory Map.

❖ HBASE-20003

# Apache HBase On Persistent Memory

❖ More memory size available - DRAM 100s of GBs. AEP even more

❖ More and More Data in Memory

o       Large memstore size and global memstore size

•       More data in memory

•       Less flushes and compactions = Less IO

•       Subsequent reads can get data from memory mostly (?)

•       More memstores size => More Java heap size => Larger GC pause issues.

•       More cell entries to CSLM.  More compares for ordering => Lower Throughput

•       Server down - more data in live WAL files for replay => Higher MTTR

o   More Java heap size –> Off heap writes using Off heap memstores

o   More cell entries to CSLM –> Compacting Memstore work by Yahoo , New faster CSLM implementation by Alibaba

o   More data in live WALs –> HBase on AEP with no WALs.  Persistent memstores LABs. Instant rebuild of memstores CSLM.
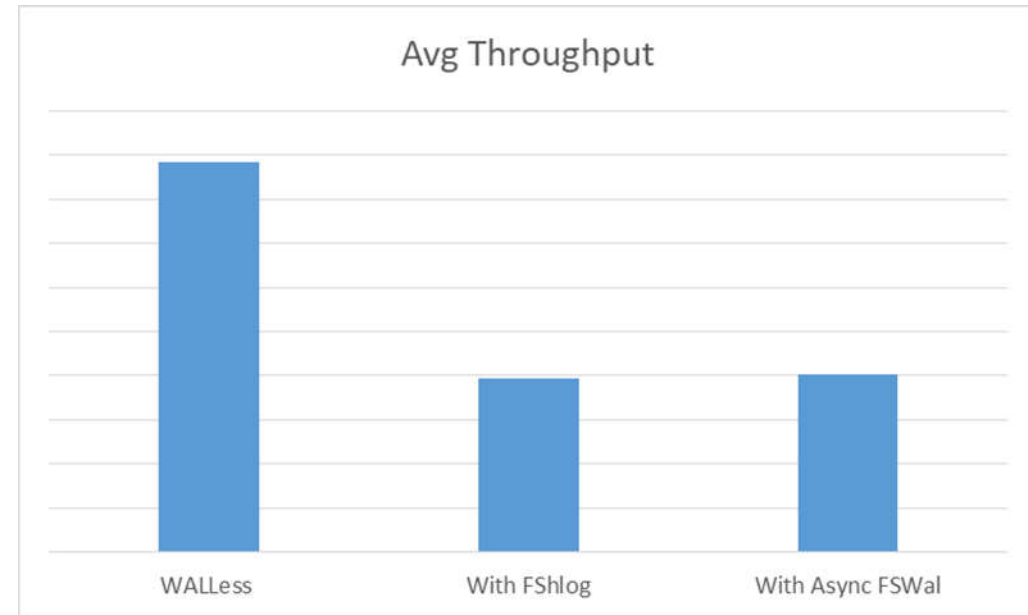
# Performance Results

❖ PerformanceEvaluation Tool

o    Write only workload

o    4 Node cluster

o    100 client thread

o    250 GB Total data

o    Single column per row

o    WAL – 3 Replicas (HDFS replicas)

o    WALLess – Primary and 2 replica regions

  Average throughput is > 2x compared to With WAL cases.
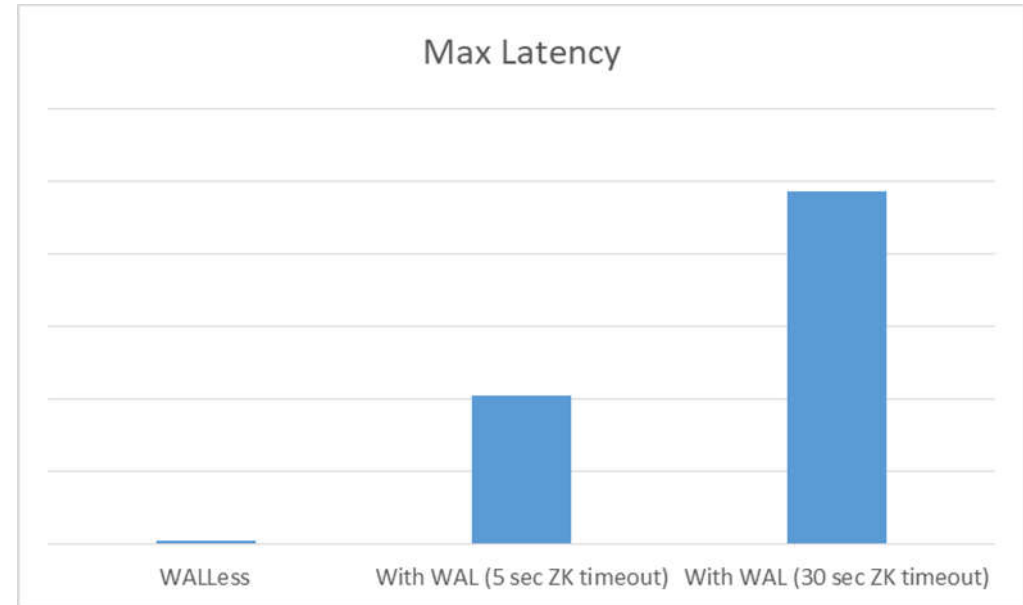
  Latency is consistent through out for the WALLess case. In WAL case the latency varies as the amount of data to 'sync' increases. (Here again with 'fsync' latency is more than 'hflush').



Avg Throughput

WALLess        With FShlog        With Async FSWal

# Performance Results

❖ PerformanceEvaluation Tool

o    Random Reads

o    4 Node cluster

o    100 client thread

o    5 secs/ 30 sec ZK session time outs

o    One RS node crash in between PE run

o

❖ Max latency is 44x larger with WAL (ZK session time out = 5 sec)

❖ Max latency is 104x larger with WAL (ZK session time out = 30 sec)



Max Latency

WALLess     With WAL (5 sec ZK timeout)     With WAL (30 sec ZK timeout)

# Apache HBase On Persistent Memory
## Project Status

❖ HBASE-20003

❖ *https://docs.google.com/document/d/1sYJS9lMZa_EMhTTOJ7y_KzVUjXdBgdKdPWmsN5p2lH0/*

❖ Write path, read path changes done in PoC

❖ Pending – WAL based features – Inter cluster replication, backup

o   Similar issue as that in HBASE-20951 (Ratis LogService backed WALs)

o   Work with this project  - HBase improvements for Cloud

❖ Testing – Full cluster restart/ Rolling restart scenarios

❖ Load balancer , AM stabilization.

o   Specially with Region replicas. Many bugs. Solving…

❖ http://pmem.io/

Thanks