

HIGH-PERFORMANCE NETWORKING

- :: USER-LEVEL NETWORKING
- :: REMOTE DIRECT MEMORY ACCESS

Hadoop

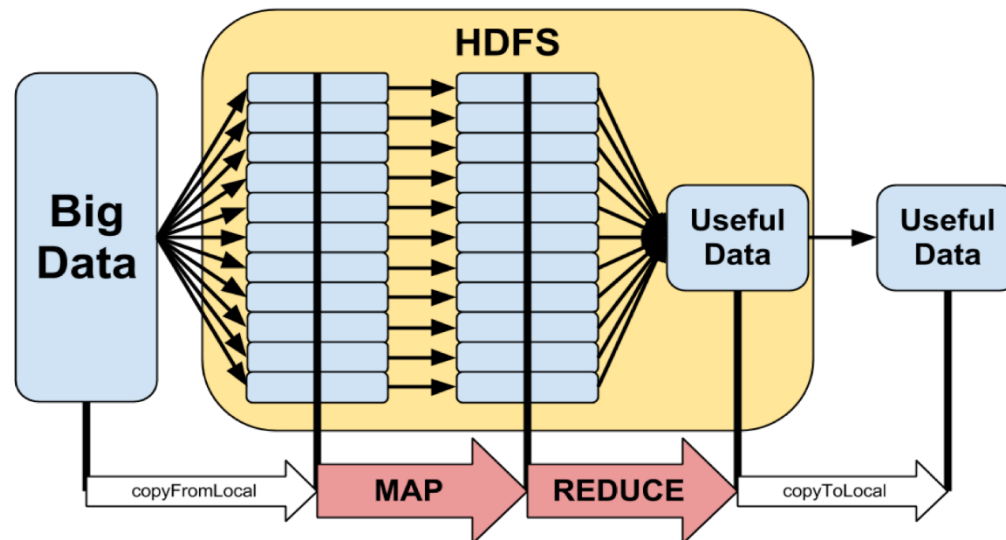
2

□ Big Data

- very common in industries
 - e.g. Facebook, Google, Amazon, ...

□ Hadoop

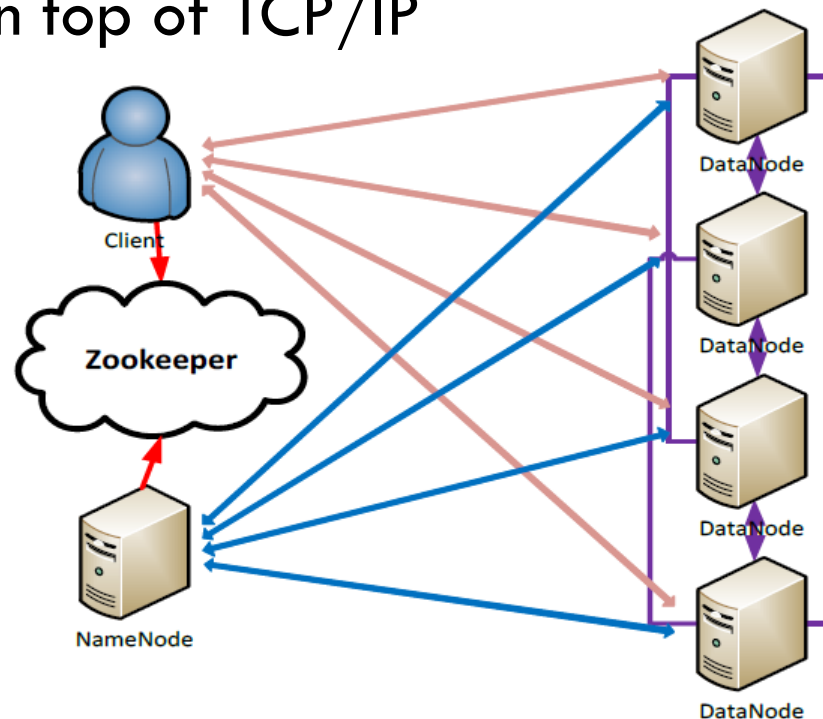
- open source MapReduce for handling large data
- require lots of data transfers



Hadoop Distributed File System

3

- primary storage for Hadoop clusters
 - ▣ both Hadoop MapReduce and HBase rely on it
- communication intensive middleware
 - ▣ layered on top of TCP/IP

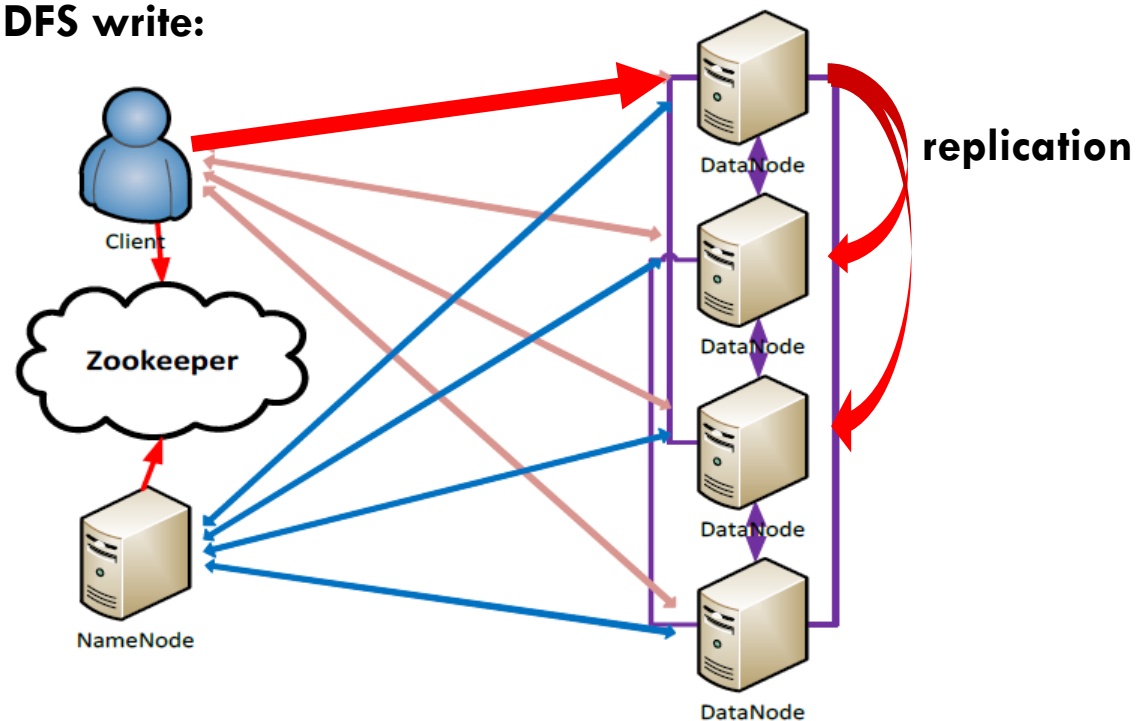


Hadoop Distributed File System

4

- highly reliable fault-tolerant replications
- in data-intensive applications, network performance becomes key component

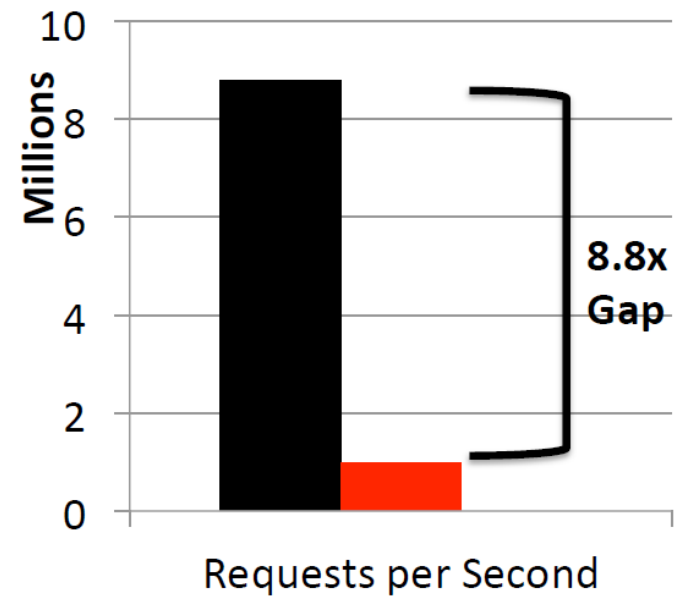
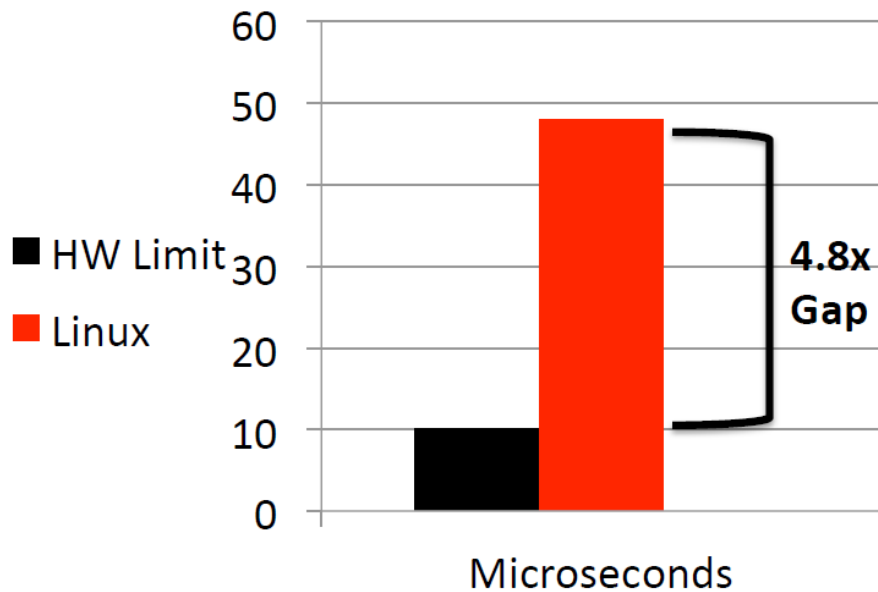
HDFS write:



Software Bottleneck

5

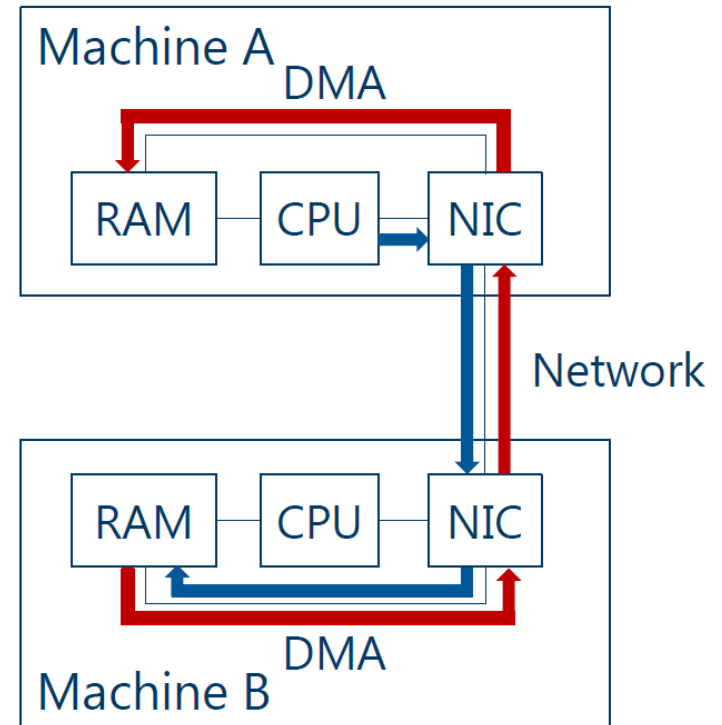
- Using TCP/IP on Linux,
 - TCP echo



RDMA for HDFS

6

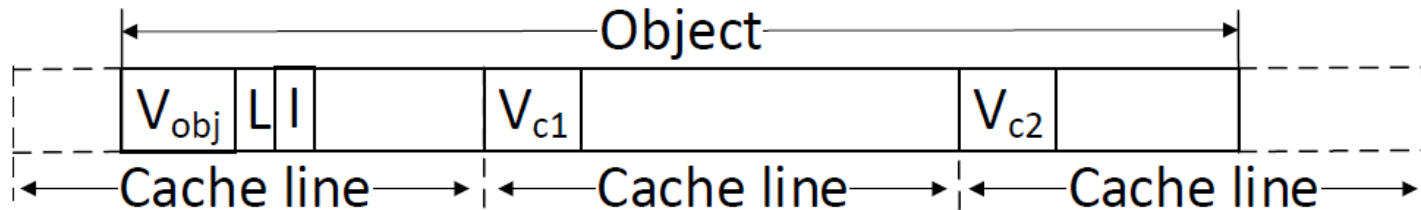
- data structure for Hadoop
 - ▣ $\langle \text{key}, \text{value} \rangle$ pairs
 - ▣ stored in data blocks of HDFS
- Both write(replication) and read can take advantage of RDMA



FaRM: Fast Remote Memory

7

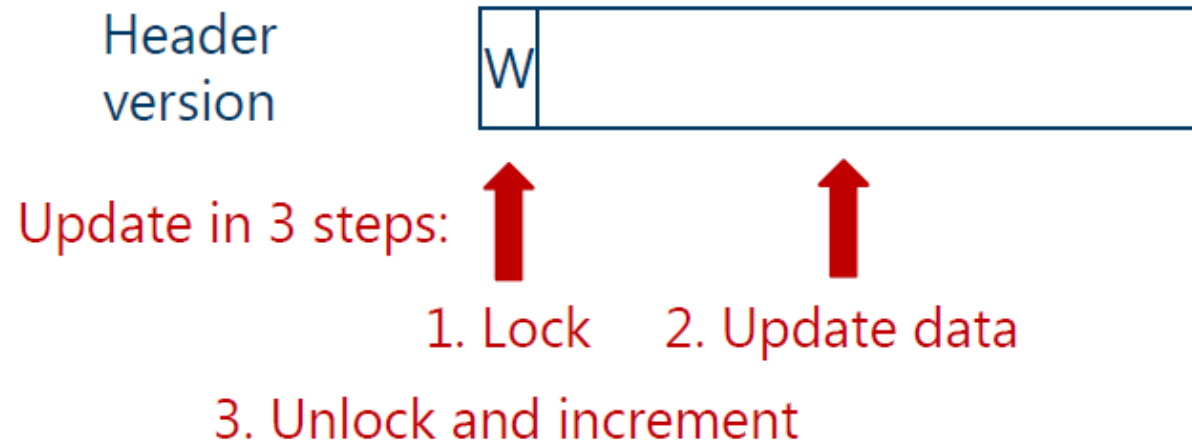
- relies on cache coherent DMA
 - ▣ object's version number
 - stored both in the first word of the object header and at the start of each cache line
 - NOT visible to the application (e.g. HDFS)



Traditional Lock-free reads

8

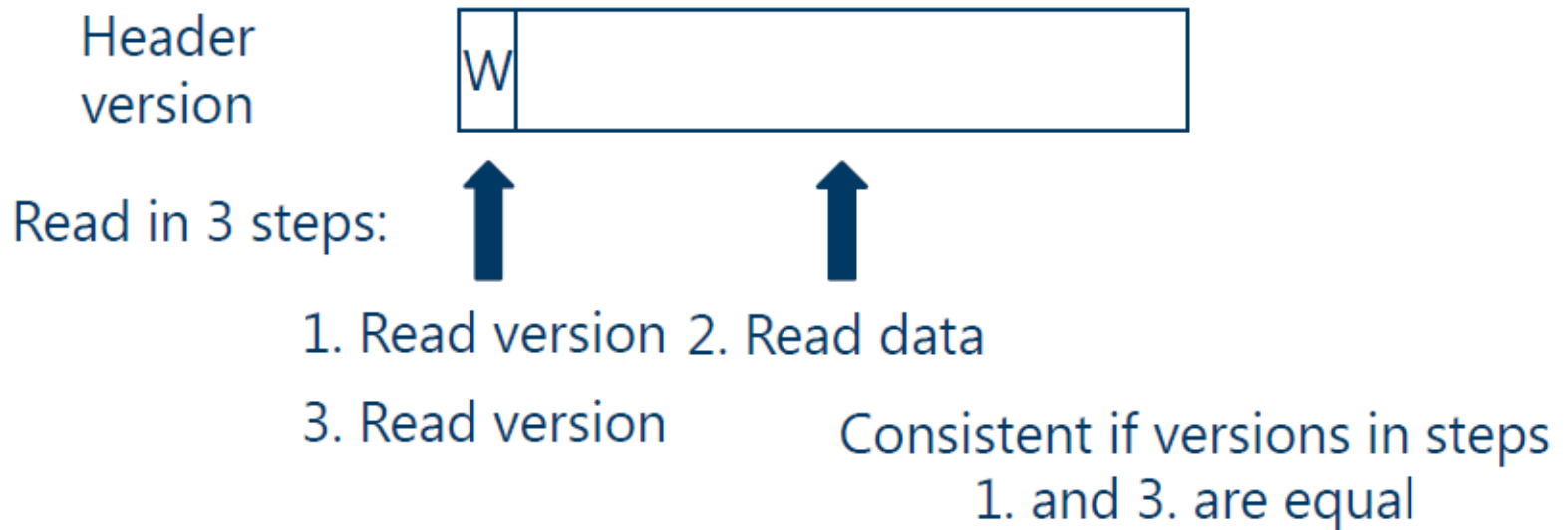
- For updating the data,



Traditional Lock-free reads

9

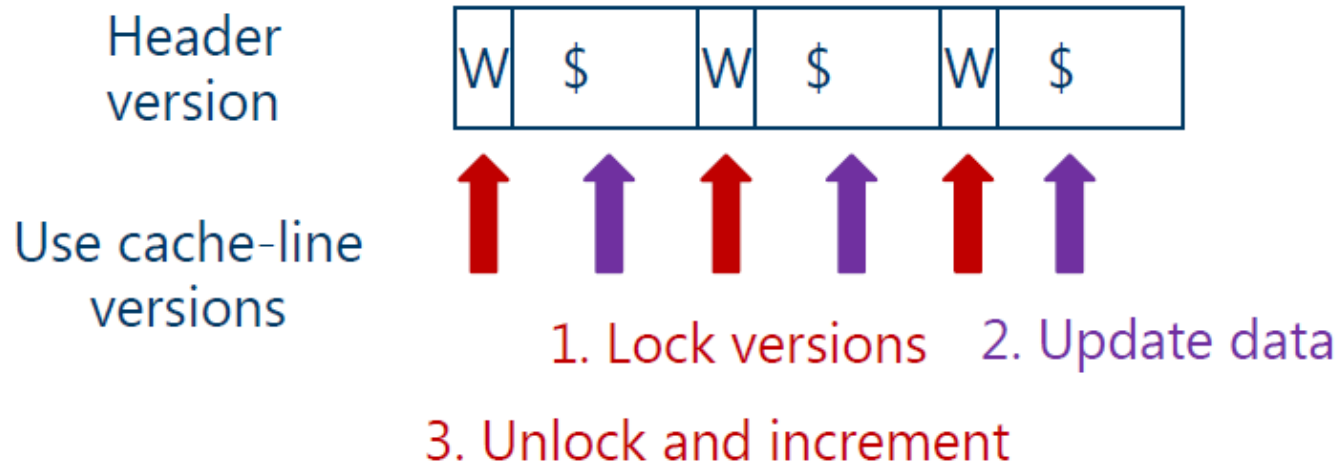
- Reading requires three accesses



FaRM Lock-free Reads

10

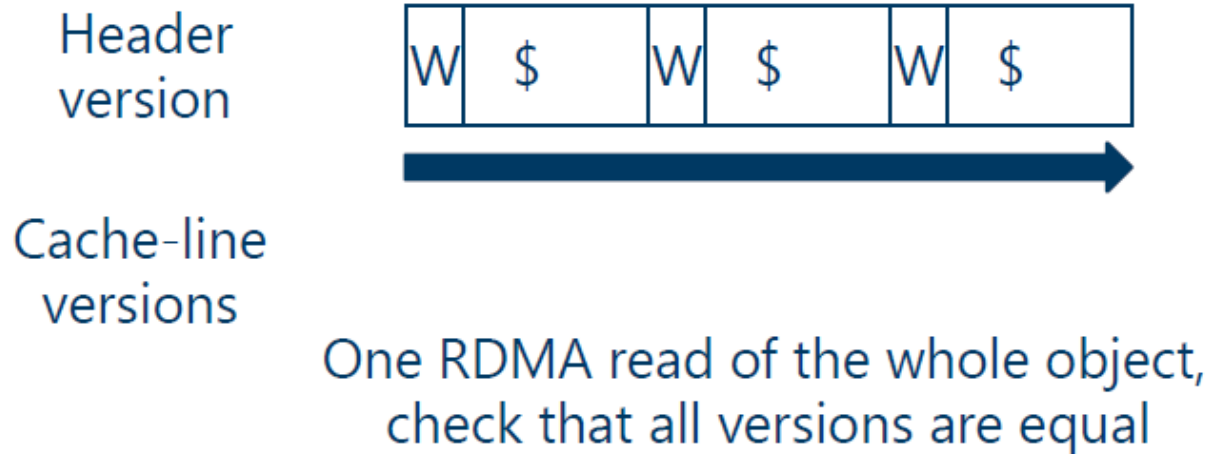
- FaRM relies on cache coherent DMA
- Version info in each of cache-lines



FaRM Lock-free Reads

11

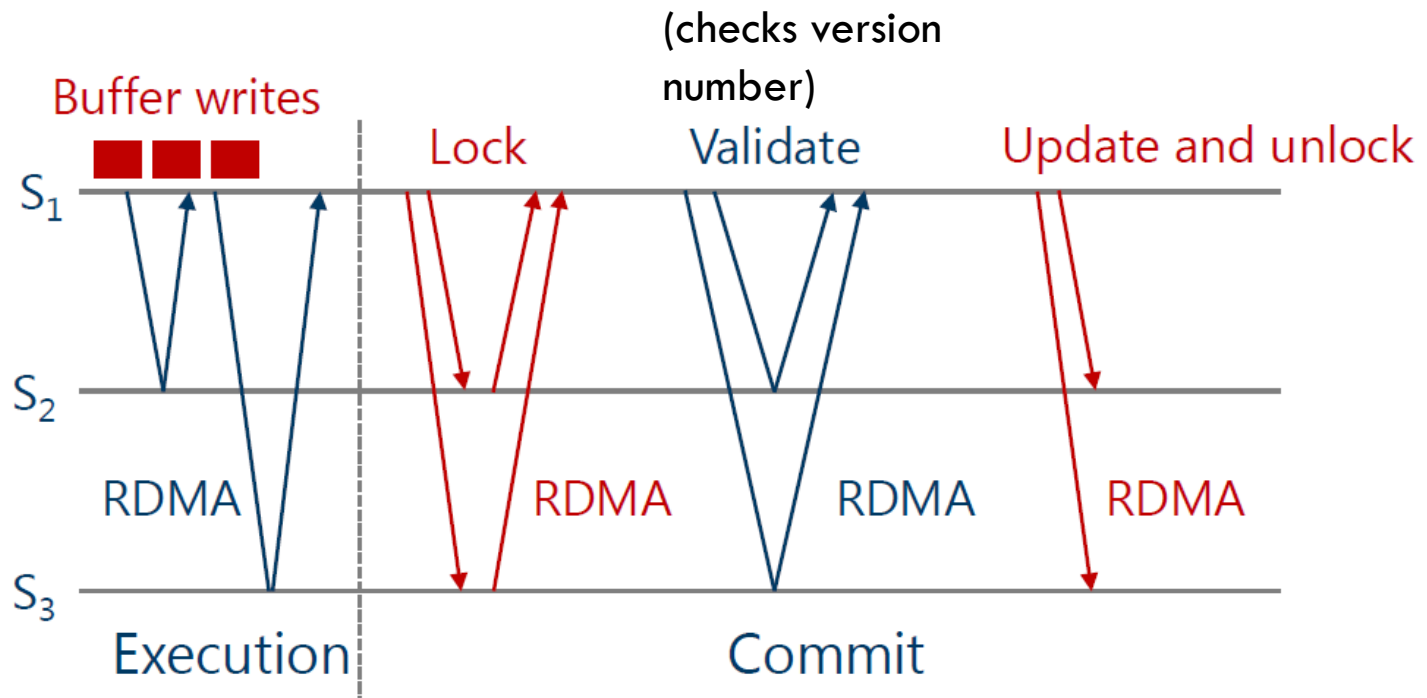
- single RDMA read



FaRM: Distributed Transactions

12

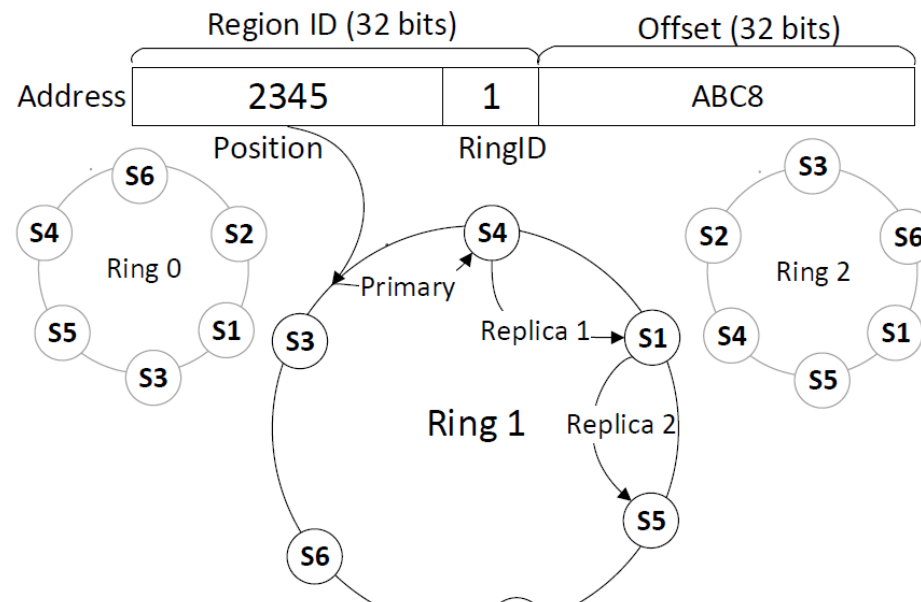
- general mechanism to ensure consistency
- Two-stage commits



Shared Address Space

13

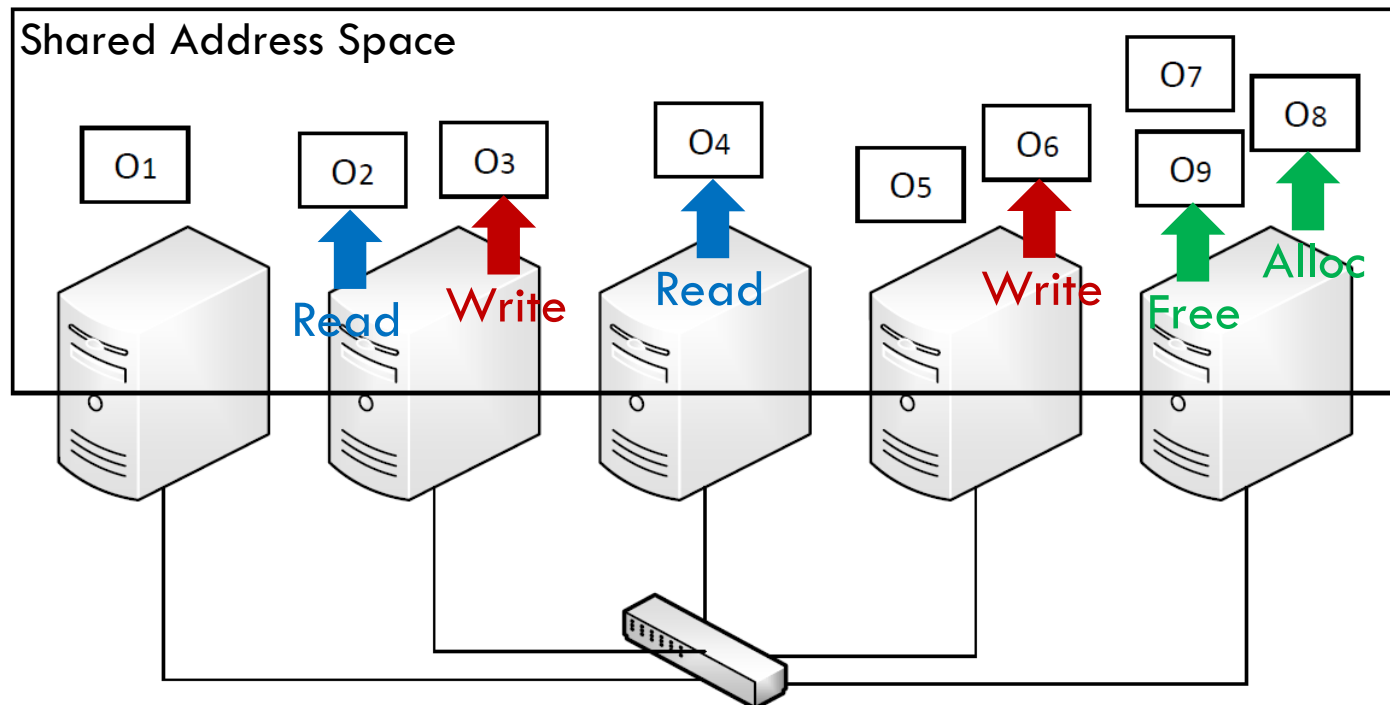
- shared address space consists of many shared memory regions
- consistent hashing for mapping region identifier to the machine that stores the object
 - ▣ each machine is mapped into k virtual rings



Transactions in Shared Address Space

14

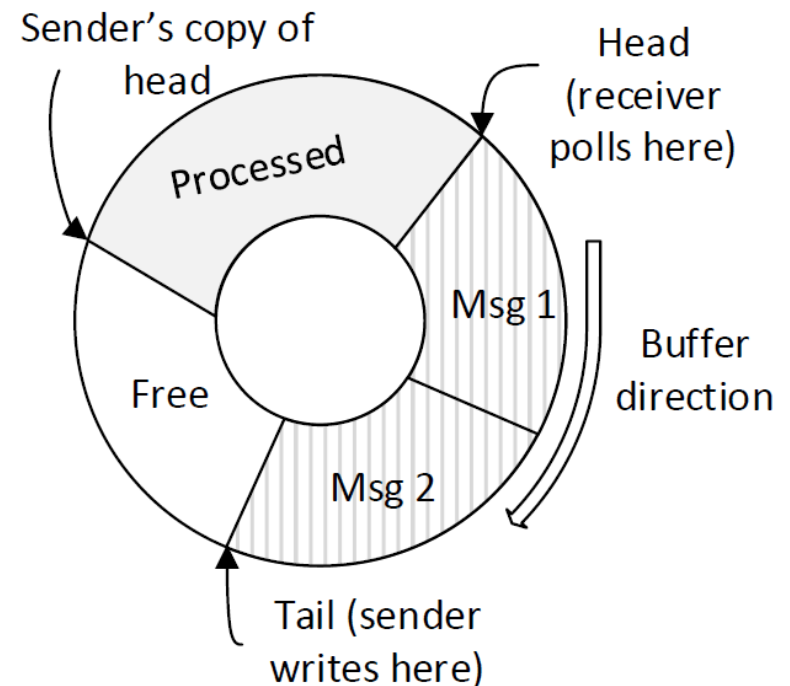
- Strong consistency
- Atomic execution of multiple operations



Communication Primitives

15

- One-sided RDMA reads
 - ▣ to access data directly
- RDMA writes
 - ▣ circular buffer is used for unidirectional channel
 - ▣ one buffer for each sender/receiver pair
 - ▣ buffer is stored on receiver



benchmark on communication primitives

16

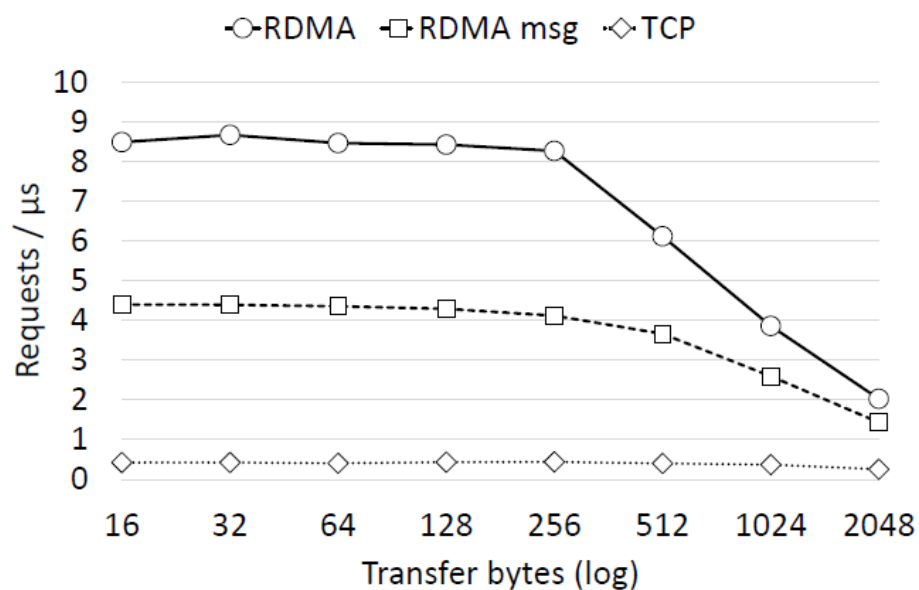


Figure 2: Random reads: request rate per machine

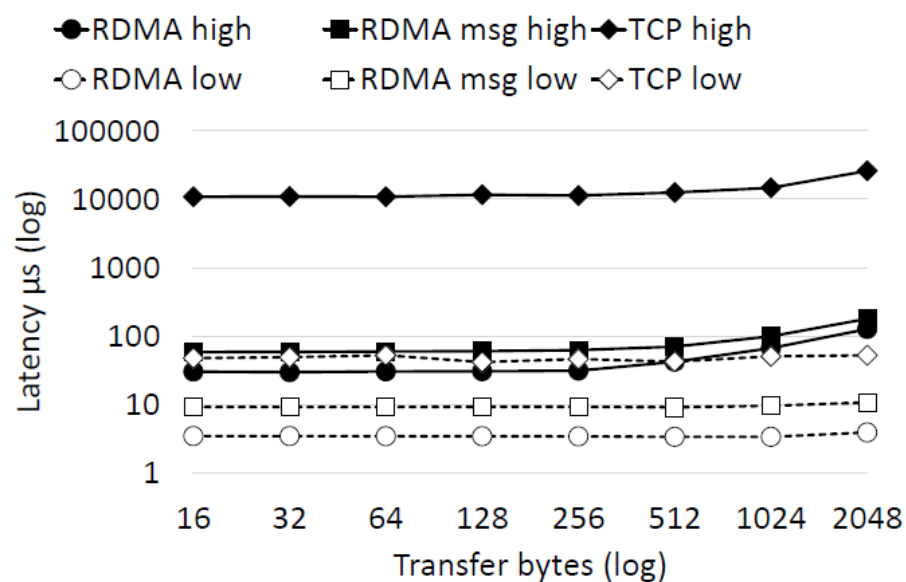
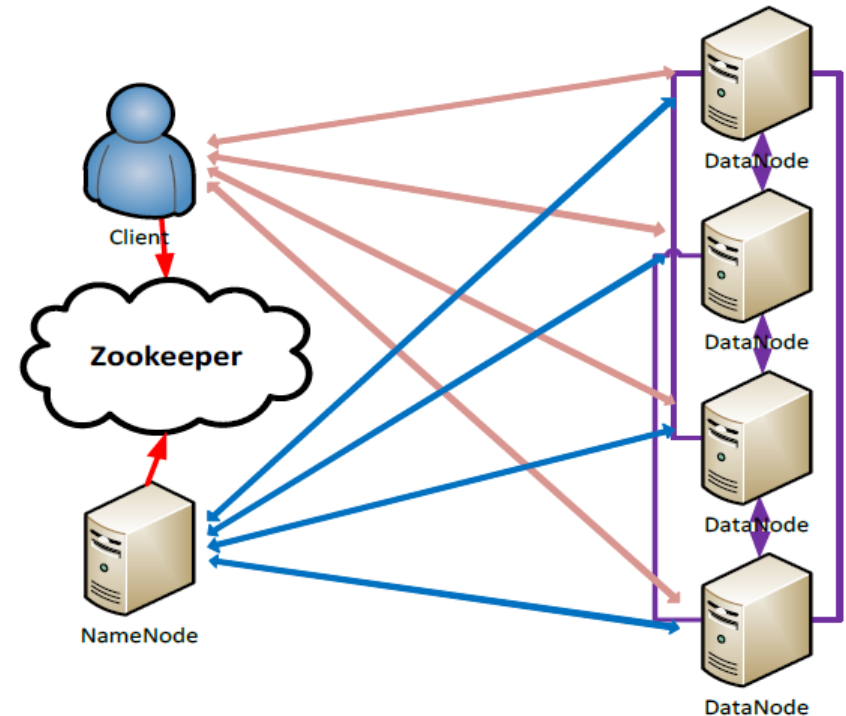


Figure 3: Random reads: latency with high and low load

Limited cache space in NIC

17

- Some Hadoop clusters can have hundreds and thousands of nodes
- Performance of RDMA can suffer as amount of memory registered increases
 - ▣ NIC will run out of space to cache all page tables



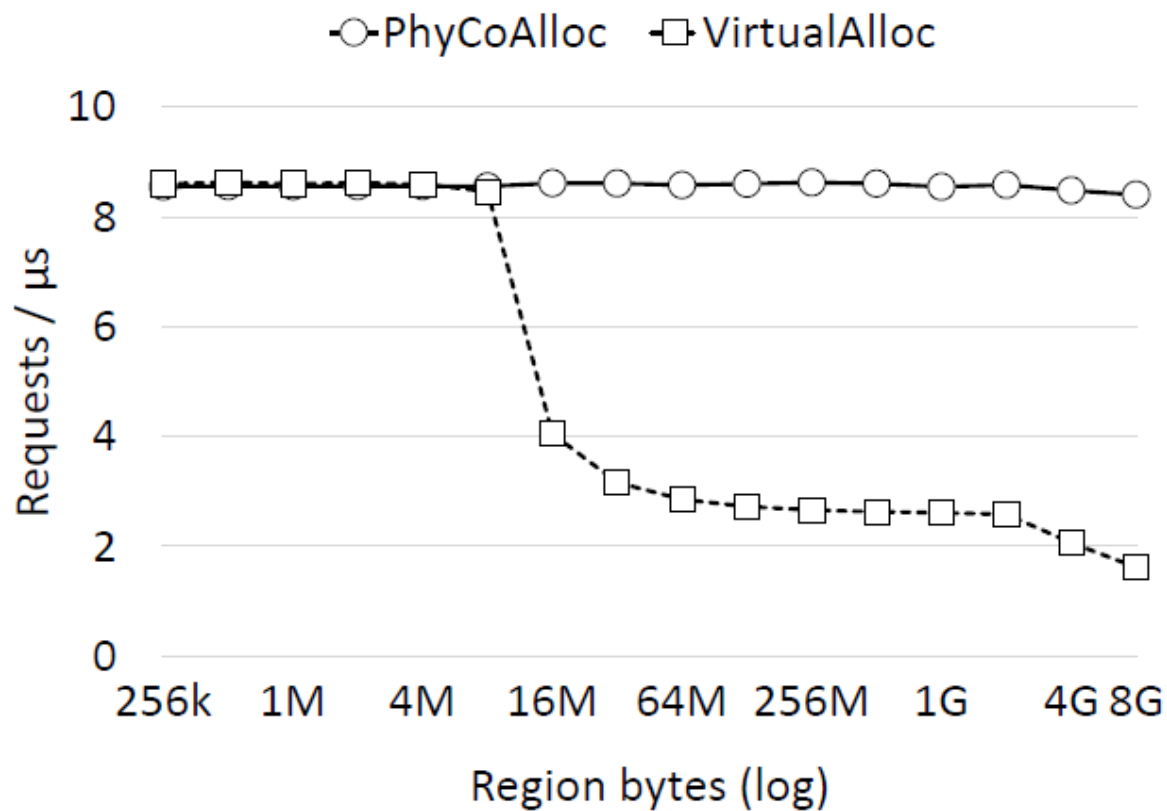
Limited cache space in NIC

18

- FaRM's solution: PhyCo
 - ▣ kernel driver that allocates a large number of physically contiguous and naturally aligned 2GB memory regions at boot time
 - ▣ maps the region into the virtual address space aligned on a 2GB boundary

PhyCo

19



Limited cache space in NIC

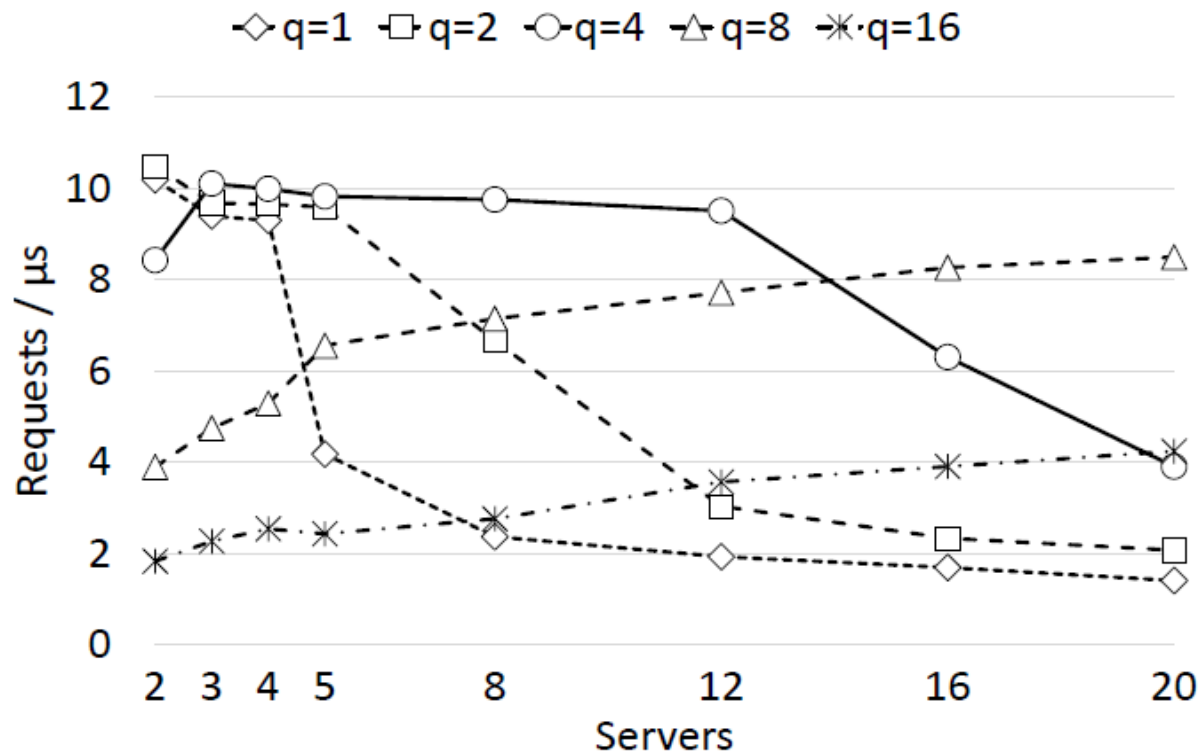
20

- PhyCo still suffered as number of clusters increased
 - ▣ because it can run out of space to cache all queue pair
 - $2 \times m \times t^2$ queue pairs per machine
 - m = number of machines, t = number of threads per machine
 - ▣ single connection between a thread and each remote machine
 - $2 \times m \times t$
 - ▣ queue pair sharing among q threads
 - $2 \times m \times t / q$

Connection Multiplexing

21

- best value of q depends on cluster size



Experiments

22

Key-value store: lookup scalability

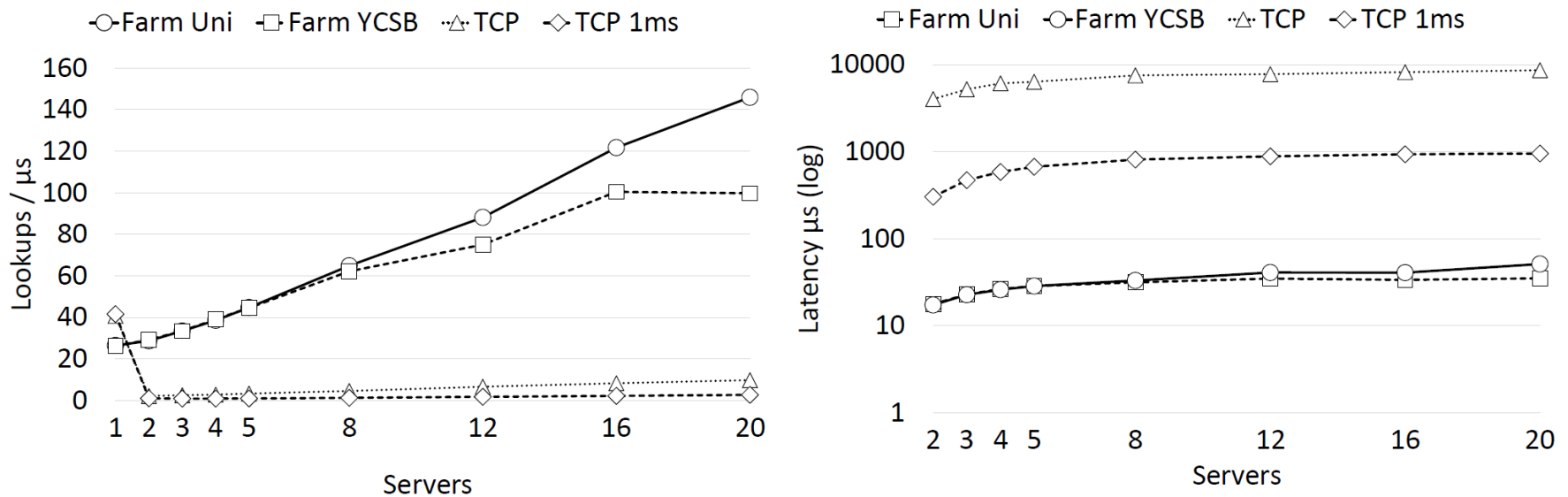
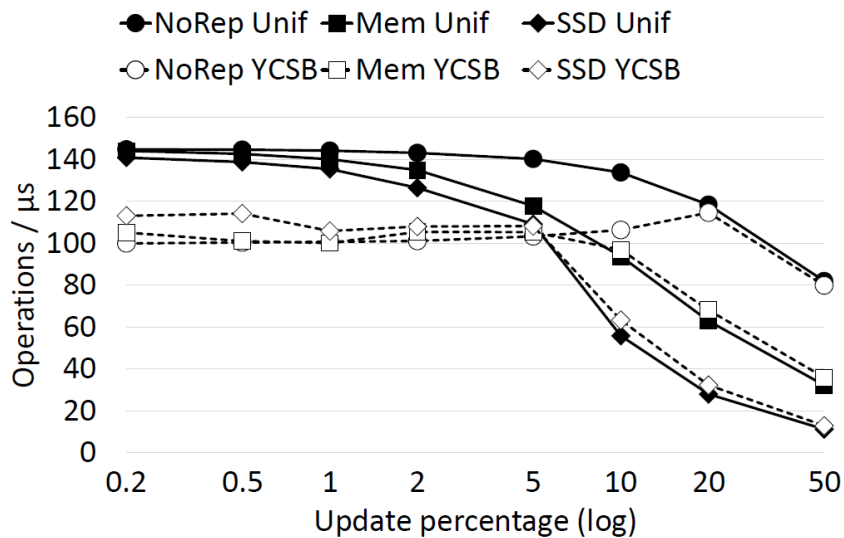


Figure 12: Key-value store: lookup scalability

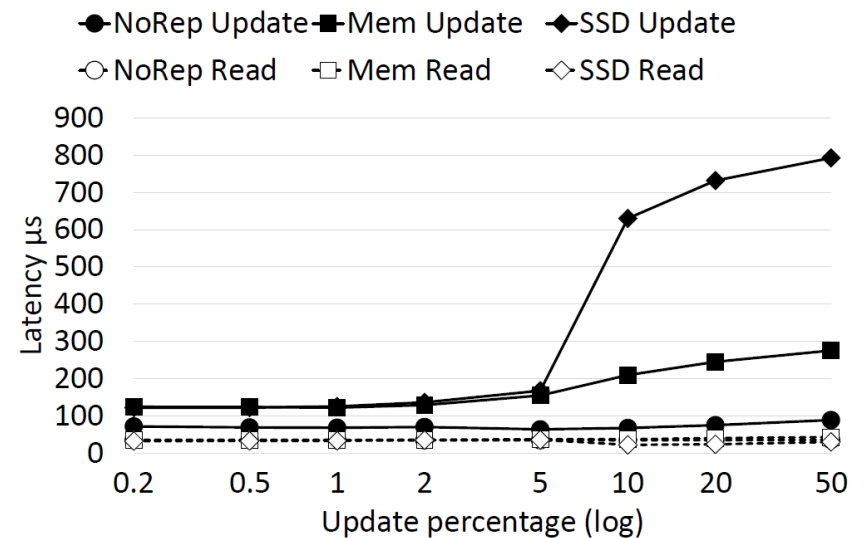
Experiments

23

□ Key-value store: varying update rates



(a) Throughput (YCSB and uniform)



(b) Latency for lookups and updates (uniform)

Figure 15: Key-value store: varying update rates