

Outback: Fast and Communication-efficient Index for Key-Value Store on Disaggregated Memory

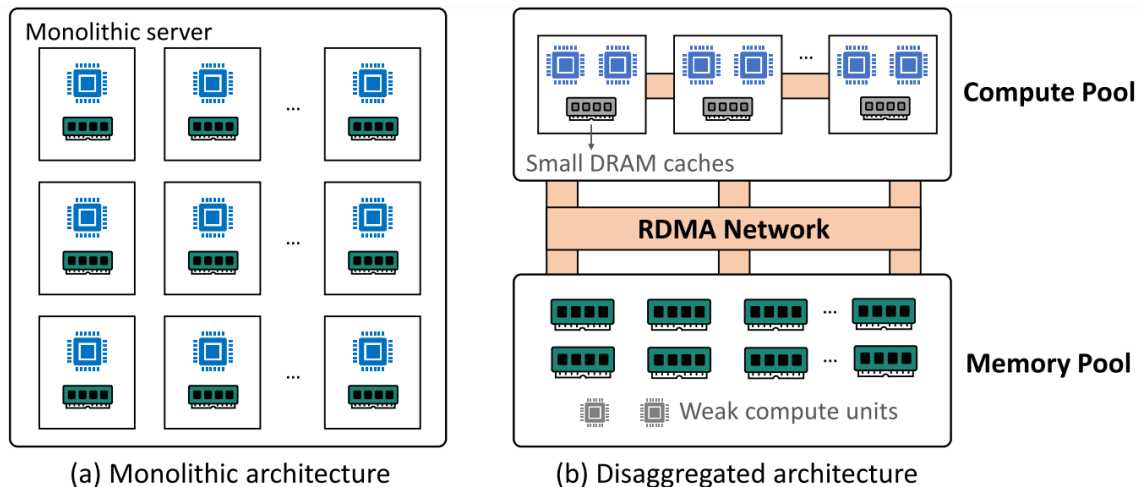
Yi Liu, Minghao Xie

Prof. Chen Qian, Prof. Heiner Litz, Prof. Yuanchao Xu

Center for Research in Systems and Storage

University of California, Santa Cruz

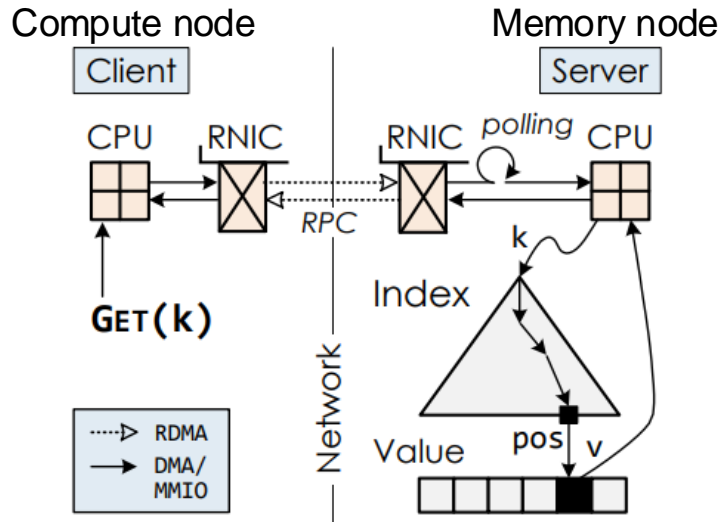
Disaggregated KV Store with Far Memory



- ❖ Disaggregated architecture decouples the compute and memory resources into independent and distributed resource pools connected by RDMA/CXL, etc.,.

[1] FORD: Fast One-sided RDMA-based Distributed Transactions for Disaggregated Persistent Memory

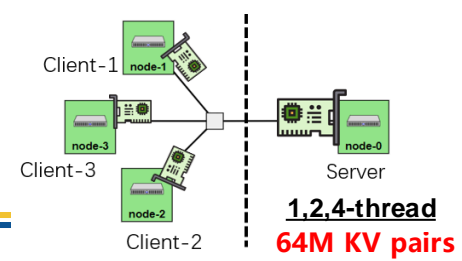
RDMA RPC based KV store



Cell (ATC'16, B tree)

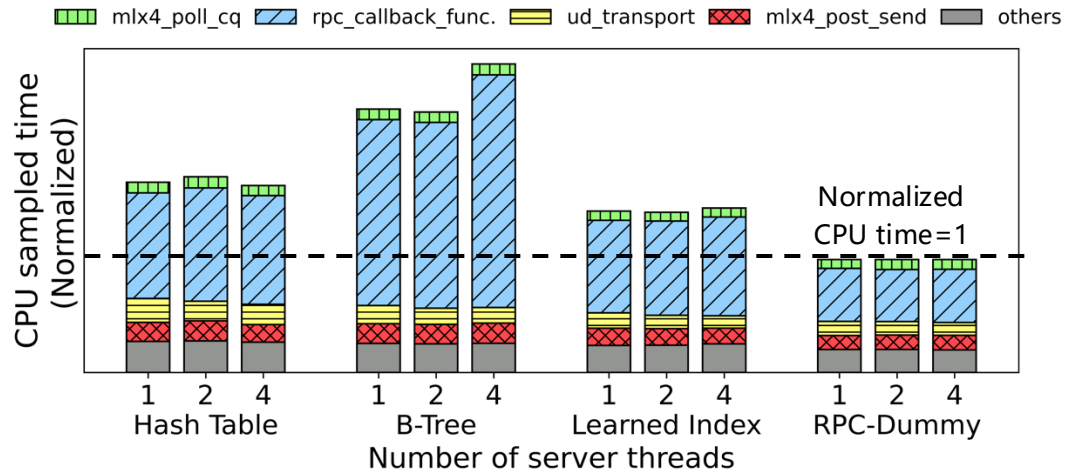
[1] Fast {RDMA-based} Ordered {Key-Value} Store using Remote Learned Cache.

Microbenchmark

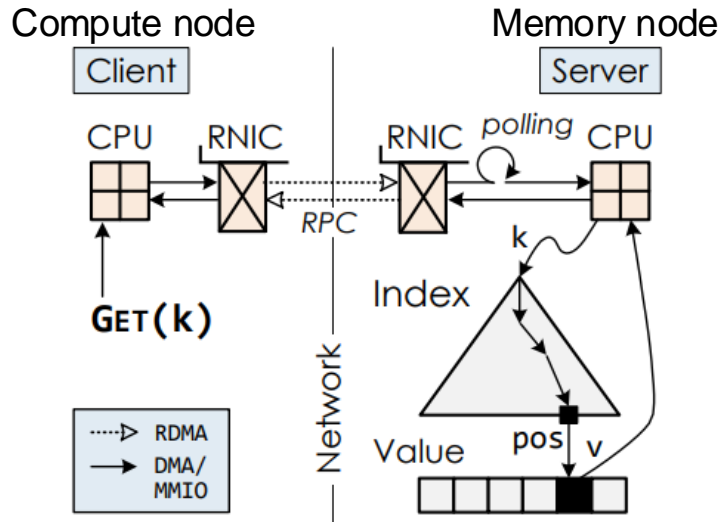


❖ Normalized CPU time for data query with:

➤ Hash table/B-tree/Learned index/Dummy



RDMA RPC based KV store



Cell (ATC'16, B tree)
FaSST (OSDI'16, Hash table)

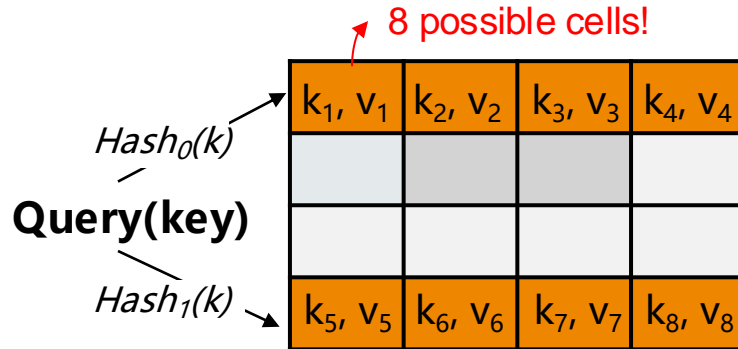
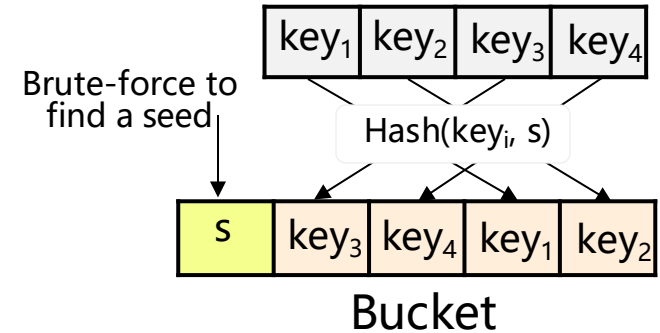


Computation on server side -> Latency

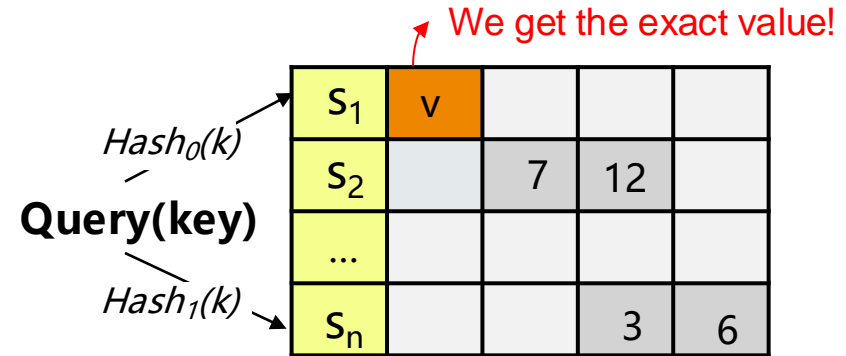
[1] Fast {RDMA-based} Ordered {Key-Value} Store using Remote Learned Cache.

Range Lookup v.s. Point Lookup

- ❖ Cuckoo Hashing / Learned Index only returns a range of possible locations
- ❖ Minimal perfect hashing (MPH) returns the exact location

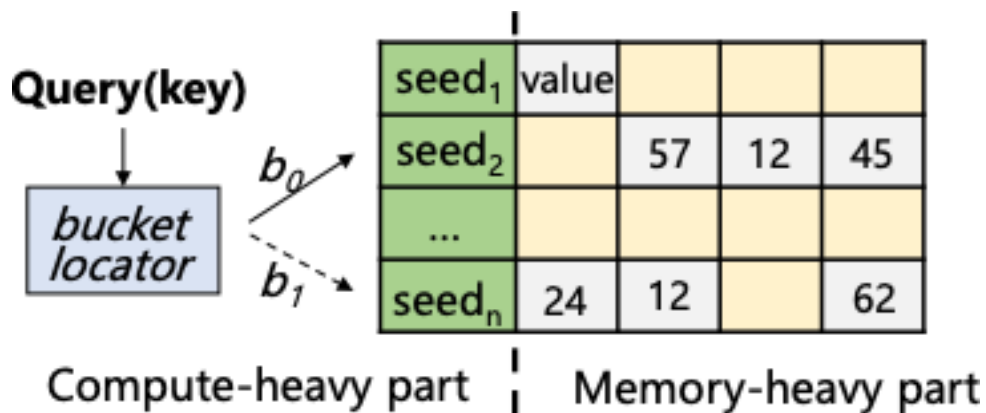


A. (2, 4)-Cuckoo Hash Table



B. MPH in Ludo Hashing

Ludo Hashing: a disaggregation fit!



Input: The Ludo lookup structure and the key k

Output: The lookup result v of k

begin

```

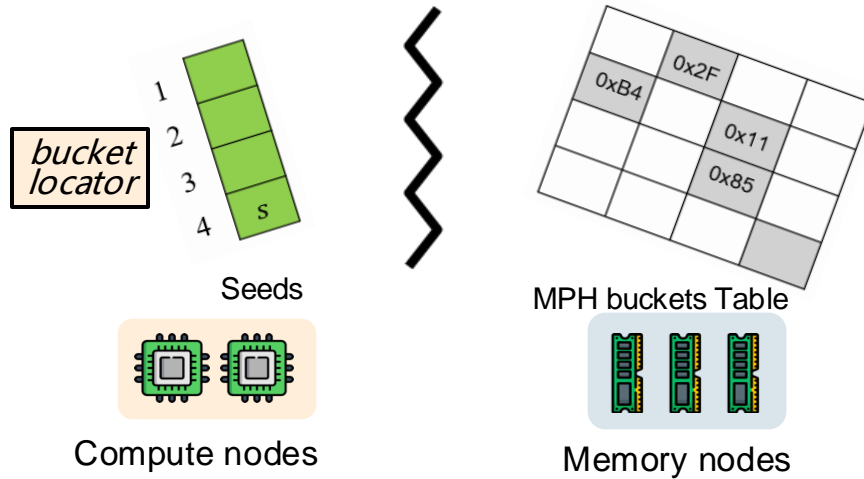
// Step I: compute bucket location
1   $b \leftarrow$  Othello lookup result of  $k$ 
2   $B \leftarrow$   $h_b(k)$ -th bucket of the table
// Step II: compute slot location
3   $s \leftarrow$  seed stored in  $B$ 
//  $\mathcal{H}_s(k) \in \{0, 1, 2, 3\}$ 
4   $v \leftarrow B.slot[\mathcal{H}_s(k)]$ 

```

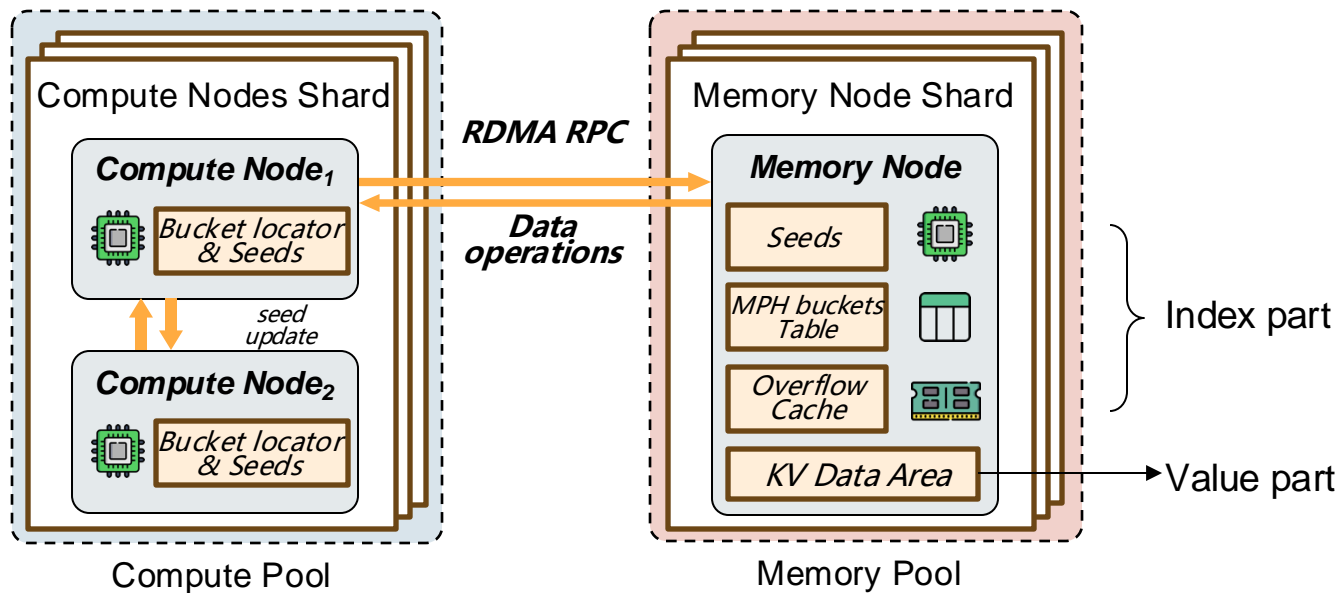
end

Ludo Hashing Lookup Algorithm

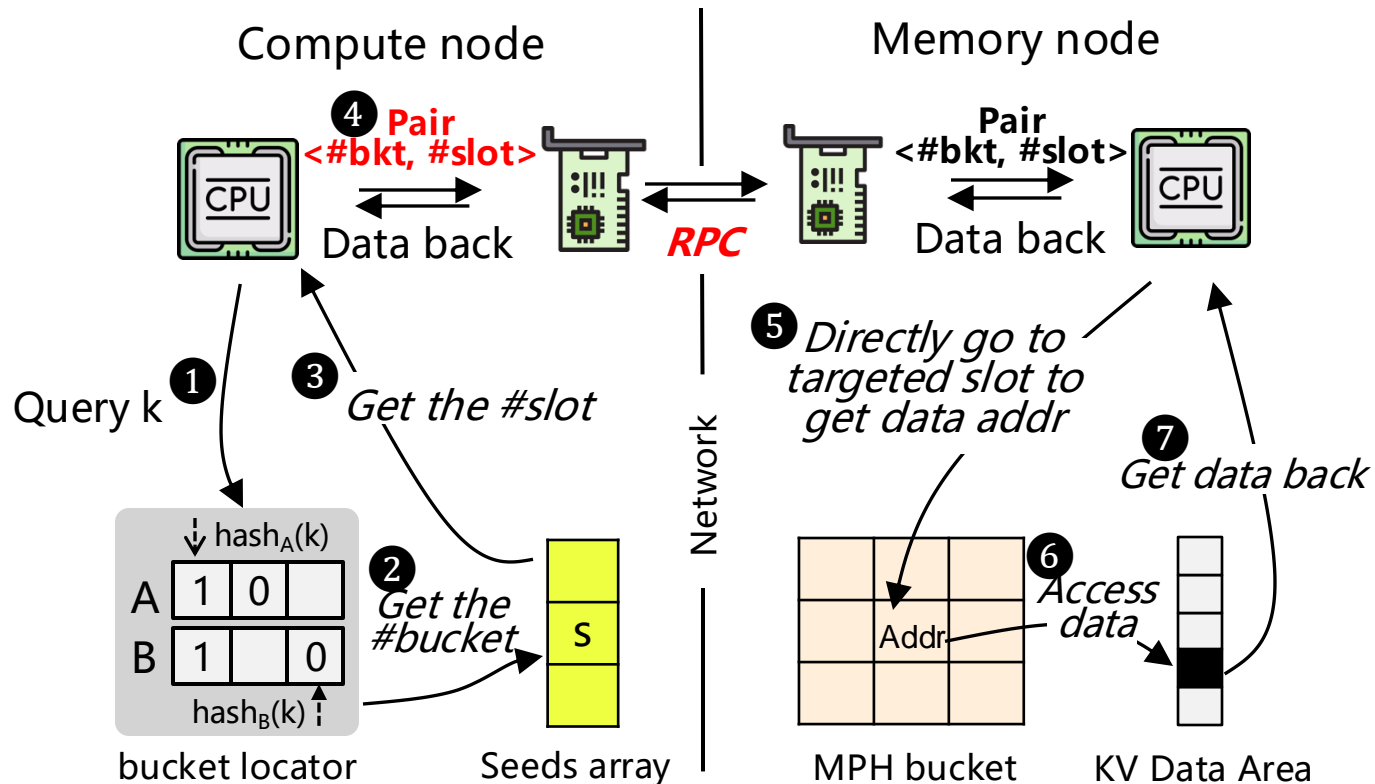
DESIGN



DESIGN

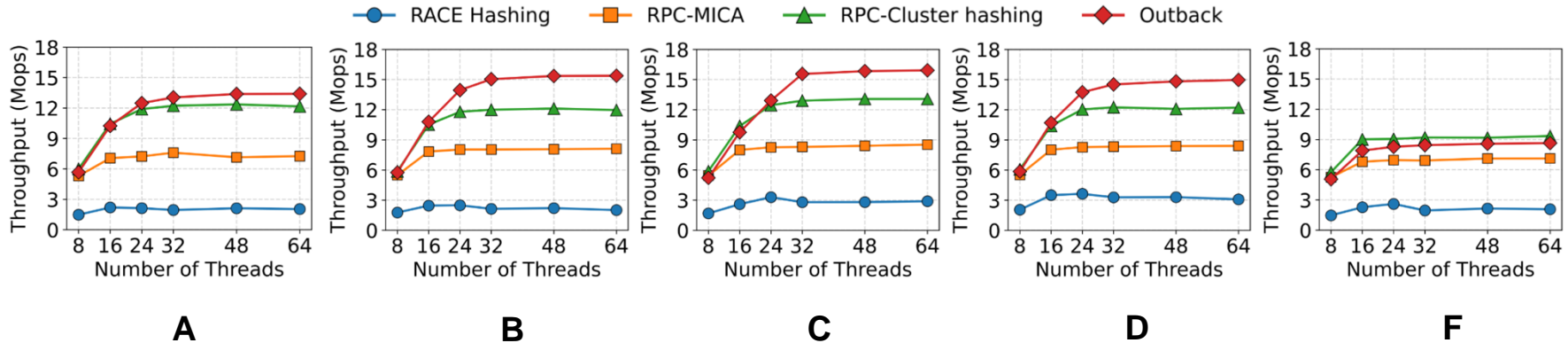


Lookup Operation



Performance Evaluation

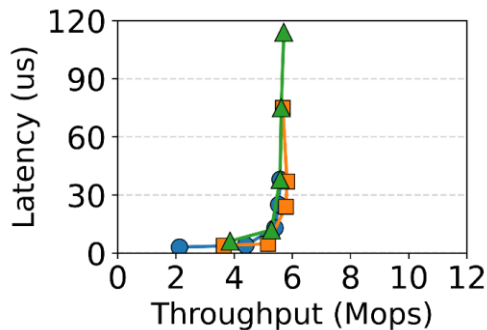
- ❖ YCSB A/B/C/D/F workloads, E is not included because this work is not optimized for data scan.
- ❖ Client threads number: 8->64, and keeps #server threads=4 and running on only 4 cores.



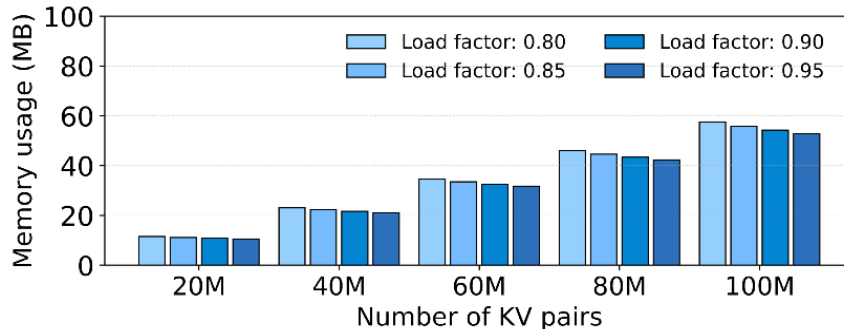
Overhead Evaluation

- ❖ Increased #Coroutines/thread does not increase the performance unlike increased #threads
- ❖ Memory cost: one-sided RDMA solutions cost 5x of MBs or more on each compute node for index caching
- ❖ Capacity impact is trivial on dataset of Facebook and Open Street Map

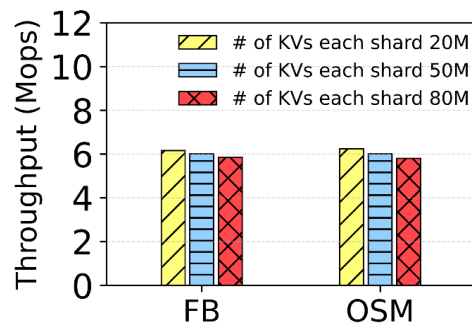
● #coroutines = 1 ■ #coroutines = 2 ▲ #coroutines = 3



Coroutine Impact



Memory Cost



Capacity Impact

Project Status



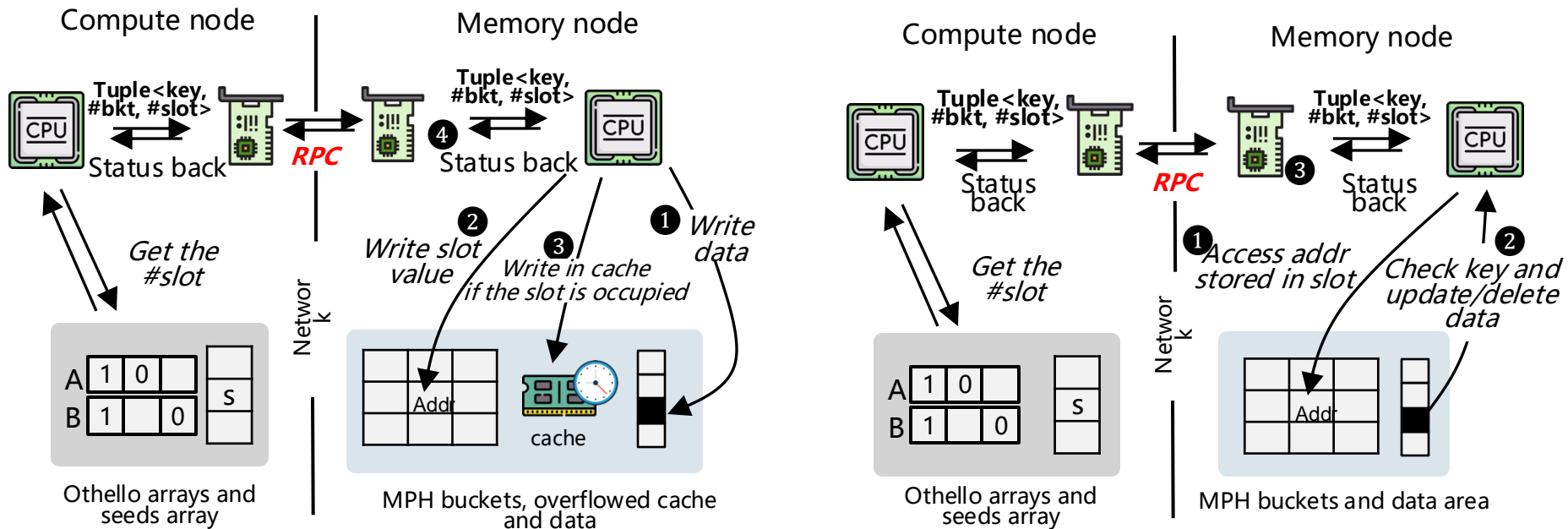
❖ Update since last IAB Meeting

- **Enhanced Experiments:** Re-ran experiments on high-performance hardware, showing Outback's scalability and throughput across varied workloads.
- **Broadened Applicability:** Demonstrated Outback's use with other data structures, supporting range queries and offloading compute tasks for devices
- **Improved Concurrency:** Implemented bucket-level locking, clarified extendible hashing resizing, ensuring efficient load balancing and shard distribution.
- Paper accepted by VLDB '25, open-sourced at <https://github.com/yliu634/outback>

Thanks for Listening

liuyi@ucsc.edu

Backup



MPH Resizing

- ❖ How to enlarge the MPH table and without stop serving data requests on the server side?

