

uKharon

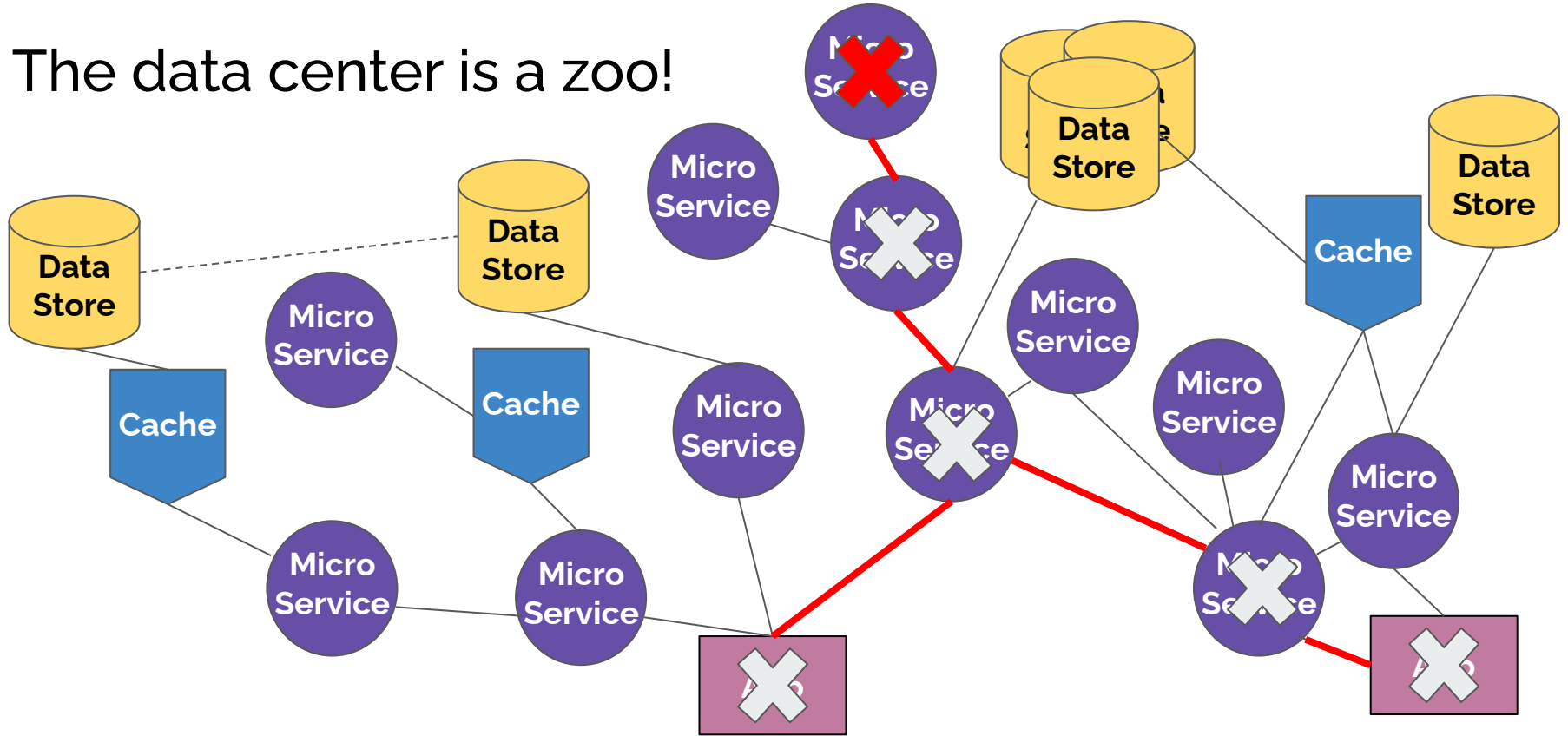
A Membership Service for Microsecond Applications

Rachid Guerraoui , Antoine Murat , Javier Picorel,
Athanasios Xygkis, Huabing Yan, Pengfei Zuo

EPFL



The data center is a zoo!



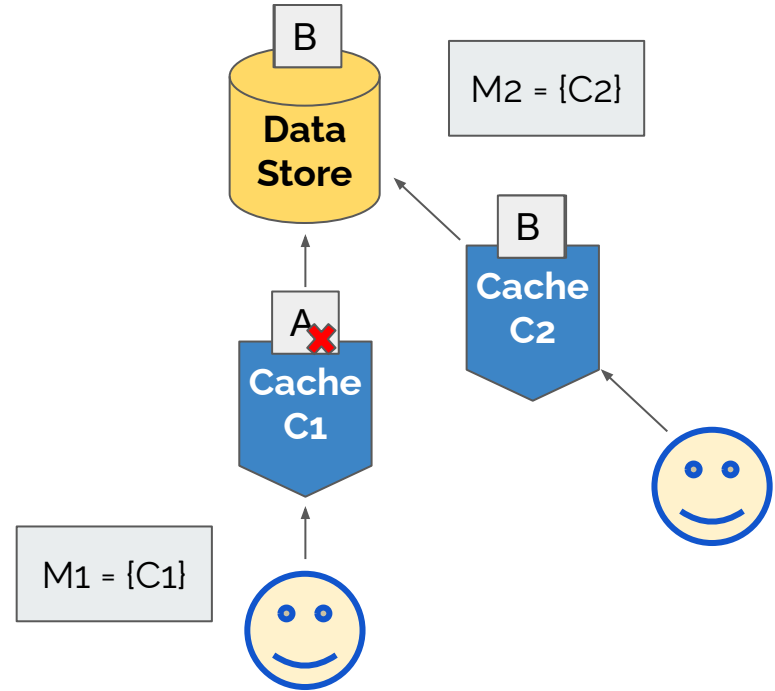
Today more than ever, failures are first class citizens!

How do we usually deal with failures?

Using etcd or ZooKeeper

Membership services:

- Are **reliable** configuration stores
- **Update** their configuration **sequence**
- **Invalidate** old memberships



The problem is NOT solved at the
microsecond scale.

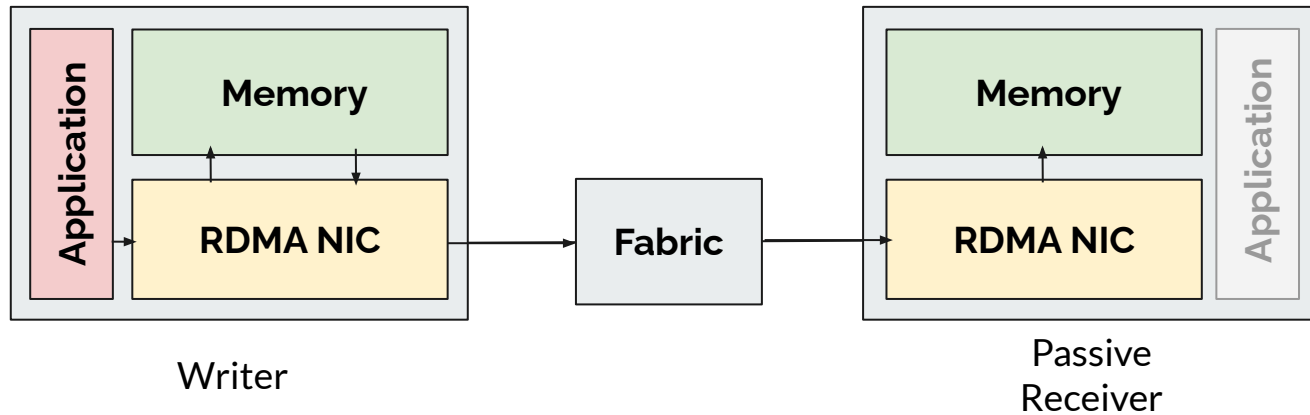
Our Contribution:

- A Microsecond-scale Membership Service
- Detects failures in **15 us**
- Updates the membership in **10us**
- Invalidates old memberships in **25us**



uKharon reacts to failures in 50us by leveraging RDMA!

Remote Direct Memory Access (RDMA)



Allows μ s-scale communication



Detects failures

Updates the membership

Invalidates old memberships

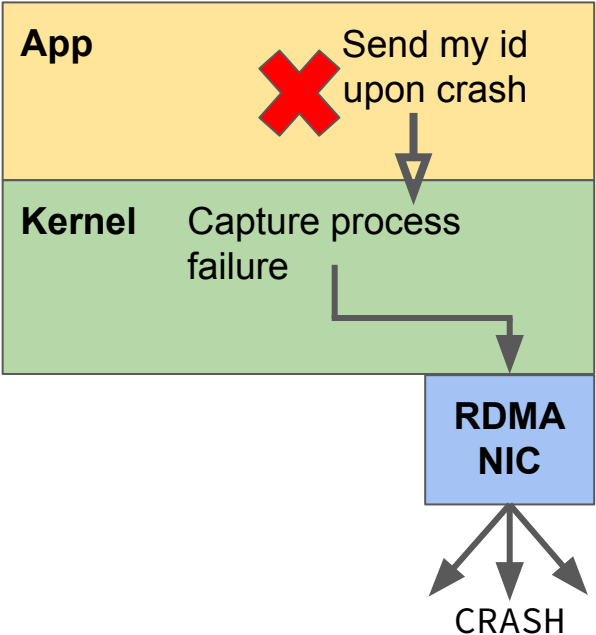
Microsecond-scale failure detection

- Timeouts do not help avoiding false positives
- Not all failures are equal

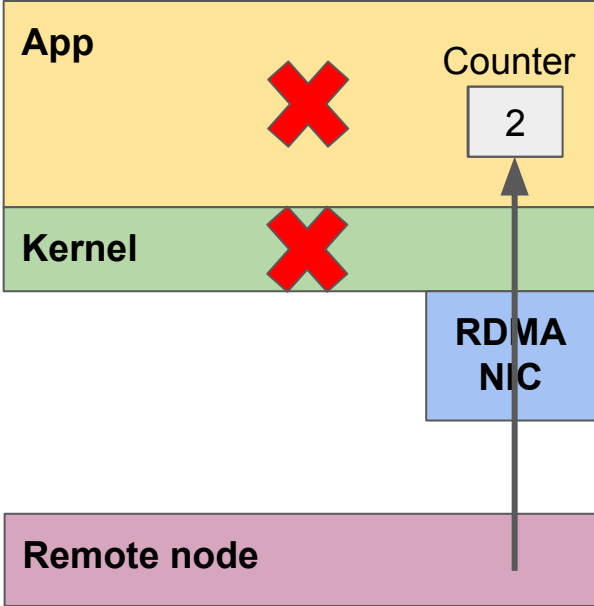
Process failures: SIGSEGV, Out-of-Memory, ...	Help from kernel → No timeouts
Kernel failures: Oops, core hang, ...	Help from NIC → No timeouts
Catastrophic failures: Power failure, NIC crash, ...	Catch-all → Timeout-based
Byzantine failures: Buffer overflow, corruption...	

Microsecond-scale failure detectors

Process failures



Kernel failures



Take network synchrony out of the equation → fast and accurate failure detection



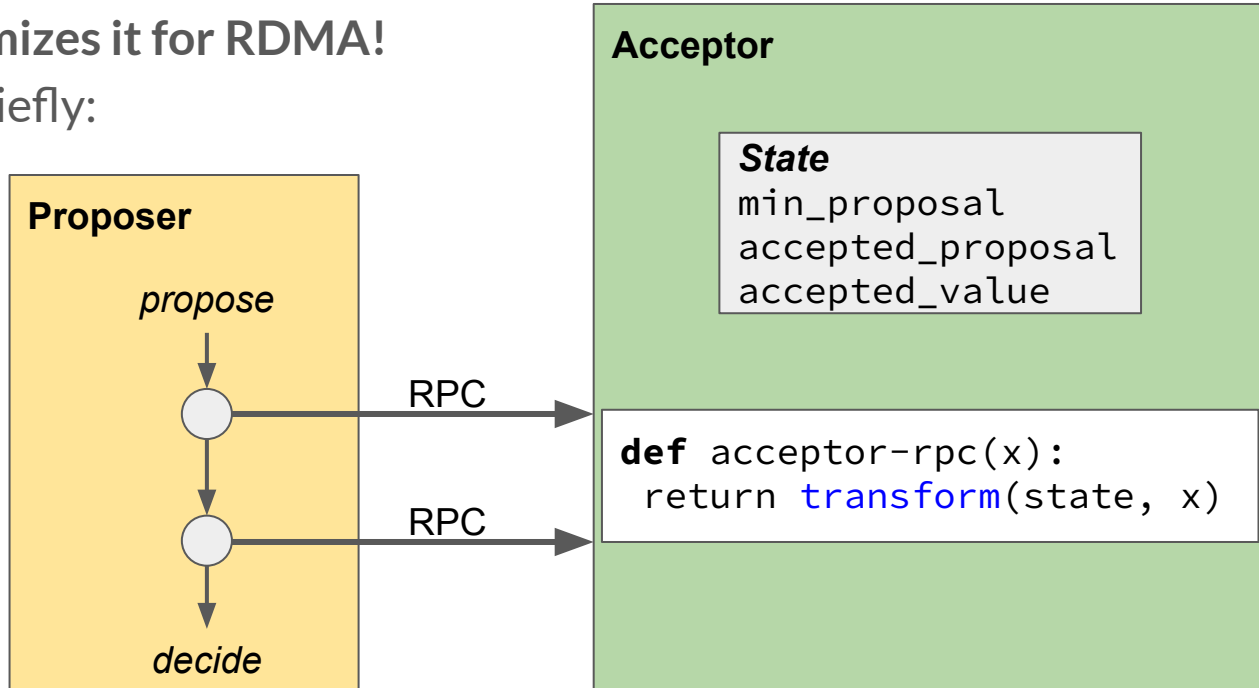
Detects failures

Updates the membership

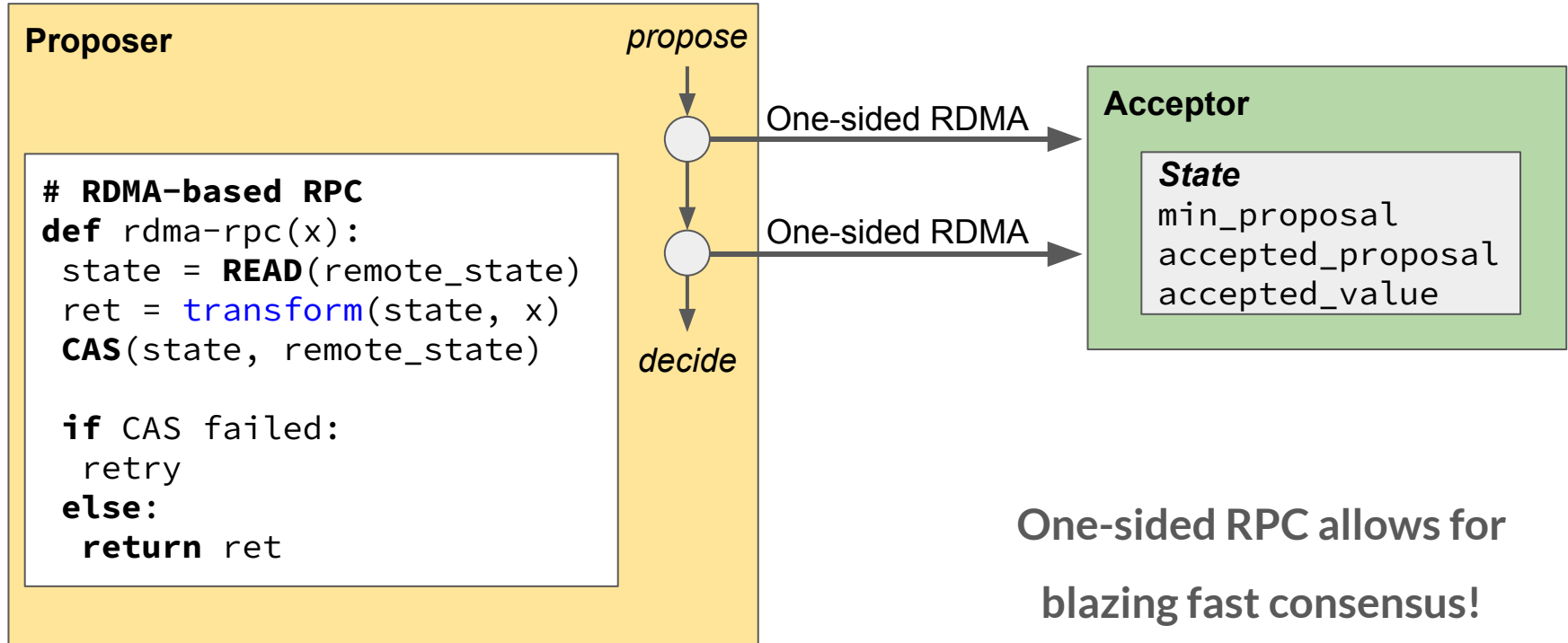
Invalidates old memberships

Microsecond-scale replication

- Uses Paxos
- But optimizes it for RDMA!
- Paxos, briefly:



One-sided Paxos





Detects failures

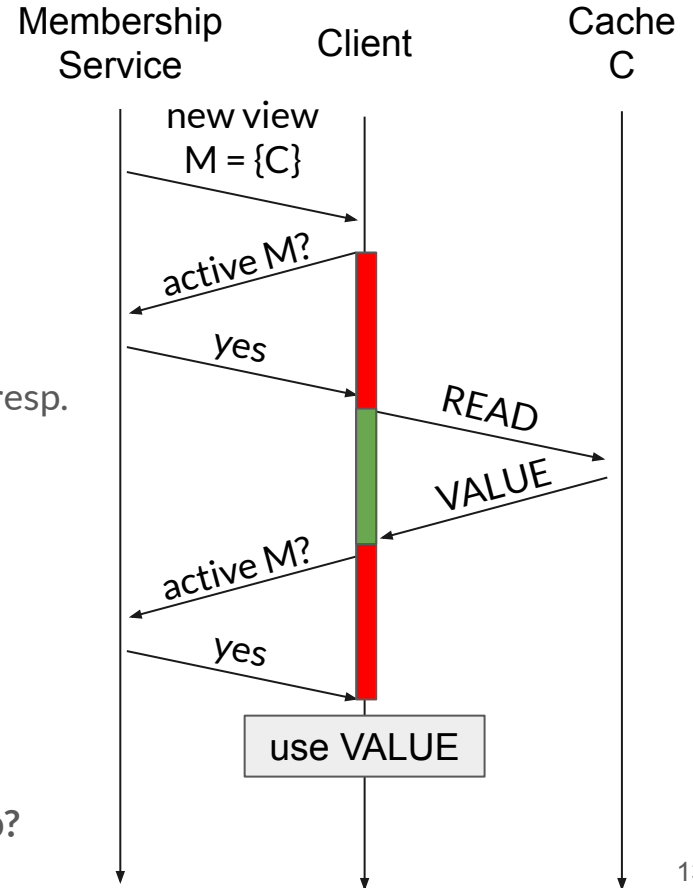
Updates the membership

Invalidates old memberships

Membership invalidation

- What is *the* active membership?
- Learn via `Active(Membership) → bool`
- `Active(M) == true → M was active between invoc. and resp.`
- High latency

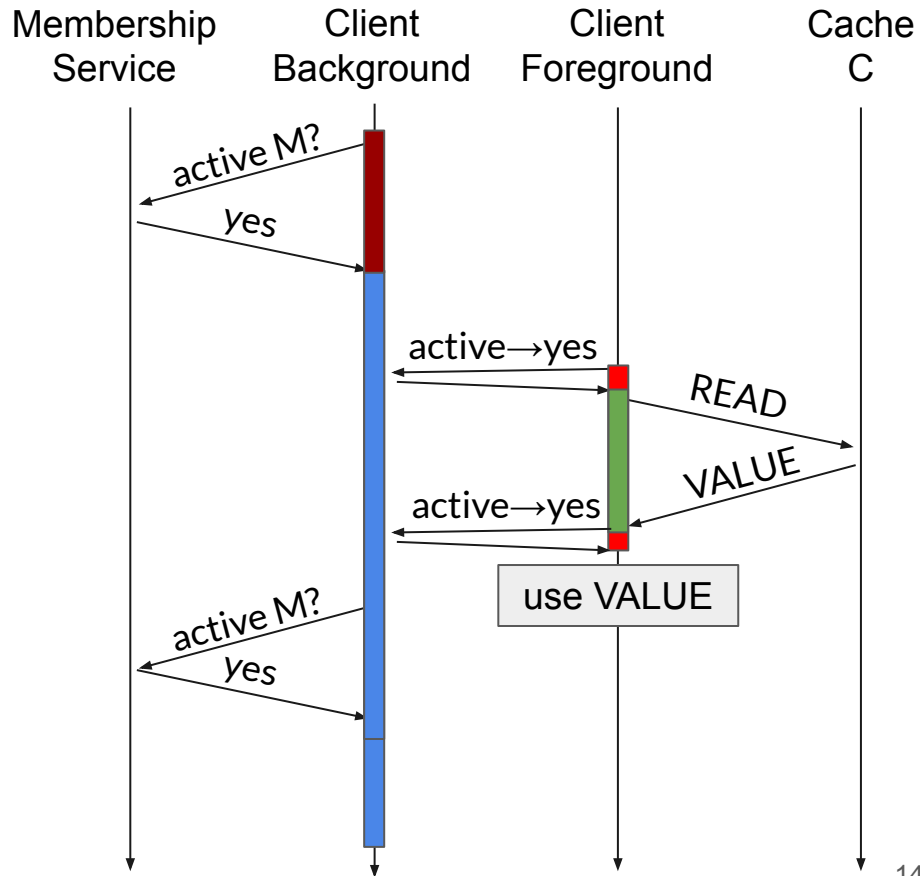
How to make **Active** cheap?



Microsecond-scale leases

- Clients lease `Active`'s response for $\sim 20\mu\text{s}$
- Renew their lease in the background
- NO synchronized clocks required
 - Only bounded clock drift for safety
- Delays view changes by no more than $\sim 20\mu\text{s}$

Leases make `Active` take $\sim 40\text{ns}$





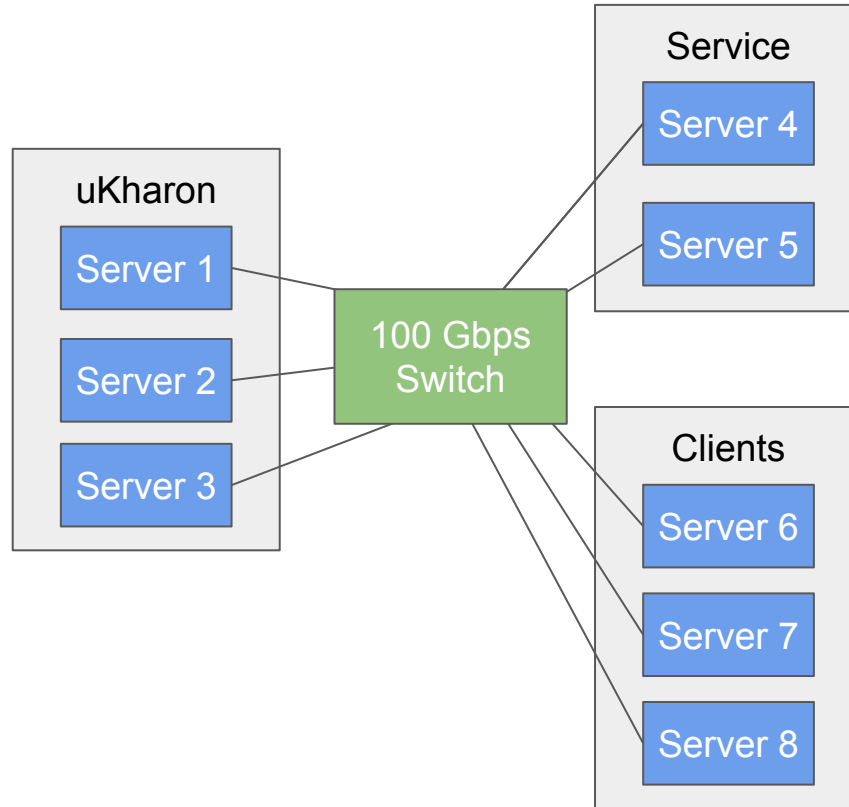
Detects failures

Updates the membership

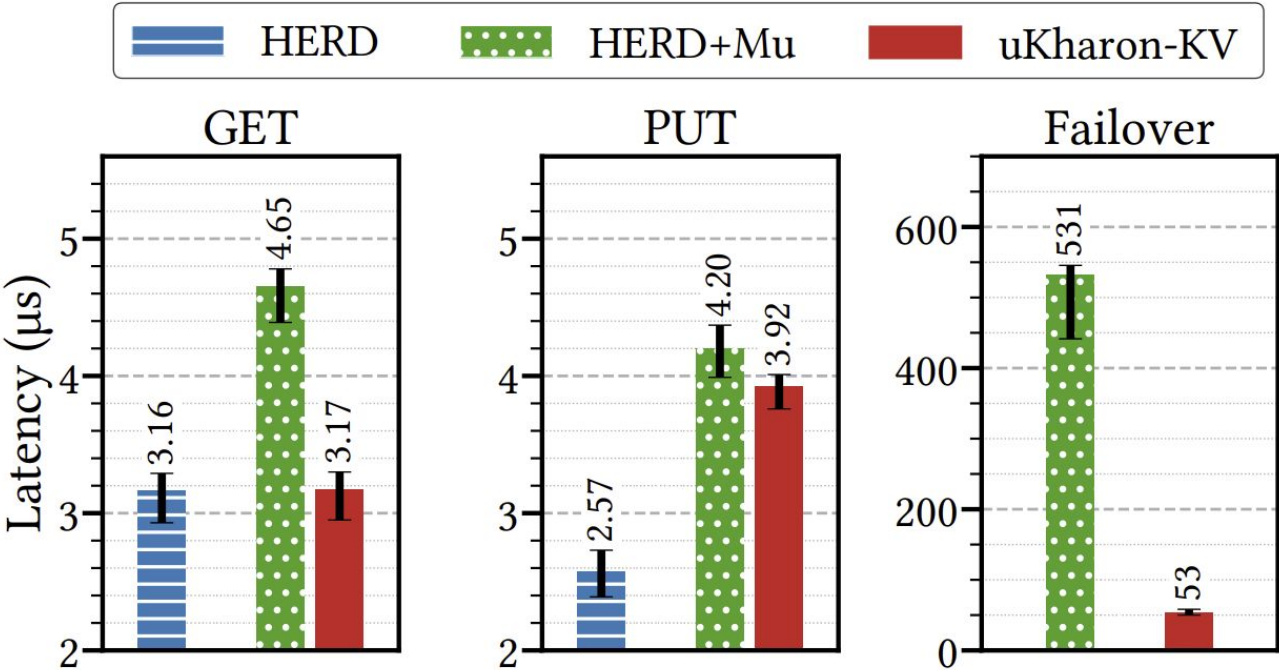
Invalidates old memberships

How does it perform?

Evaluation: setup

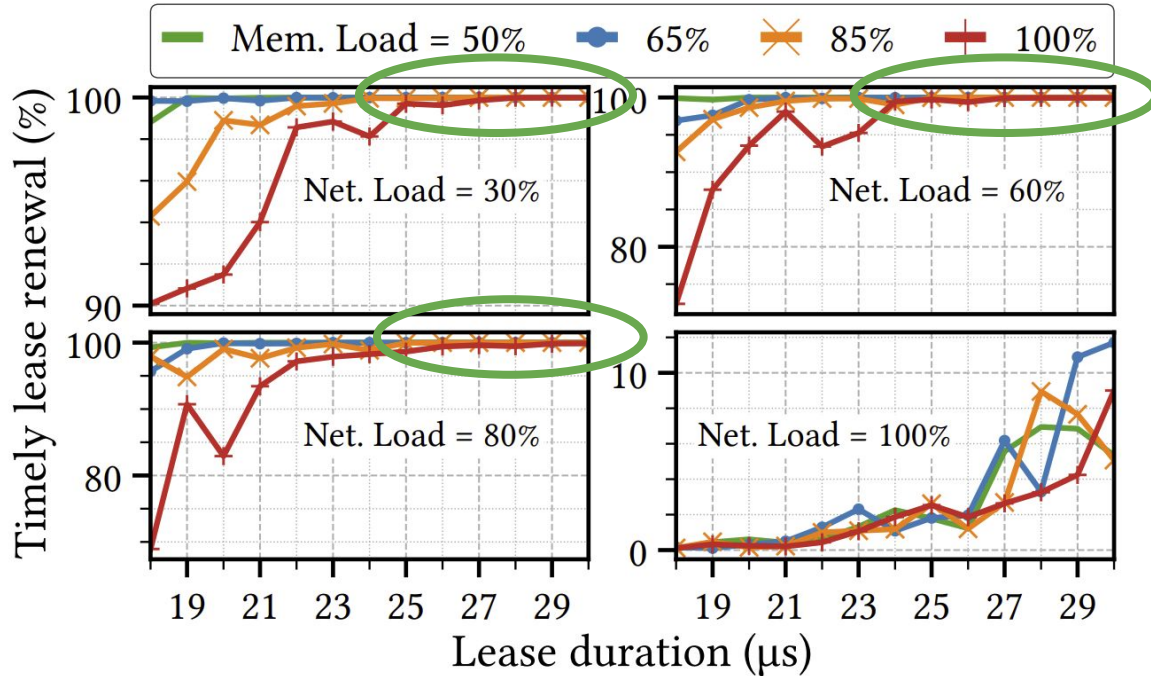


Evaluation: Replicating a KV-Store



uKharon helps beating the state of the art!

Evaluation: Are leases renewed in time?



Microsecond leases are stable!

Conclusion

- uKharon:
 - A membership service for μ s-fast failover (down to 50 μ s)
 - Easy to integrate
 - Only 40ns latency overhead



github.com/LPD-EPFL/ukharon

Check out our paper for more details!