# Beyond reCAP: Local Reads and Linearizable Asynchronous Replication

A. Katsarakis*†, E. Gioratmis*♣, V. Gavrielatos†, P. Bhatotia♣, A. Dragojevic♦, B. Grot♠, V. Nagarajan♠, P . Fatourou♥

† Huawei Research, ♣TU Munich, ♦Citadel Securities, ♠University of Edinburgh, ♥University of Crete and FORTH, *Equal contribution

## Motivation

**Online Services & Cloud Applications**

Characterized by
- Many **concurrent requests**
- **Read intensive** workloads
- Need for **data reliability**
  → run on fault-prone h/w

**Fault-tolerant Replicated Datastores**

- **Crash-tolerance**: data are replicated
- **High performance**: especially for reads
- **Strong consistency** under **asynchrony**
  → correct — even if timeouts do not hold

**Crash-tolerant Replication Protocols**
determine actions for *reads* and *writes*

**Ideal features**
1. **Linearizable**
2. **Asynchronous**
3. **Local reads**: for max perf.

## Theory



**Crash-tolerant protocols: 2 out of 3**

Linearizable

Asynchronous — Local reads

**RA protocols:**
- **Remote (costly) reads**
+ **Linearizable**
+ **Asynchronous**

**LS protocols:**
- **Synchronous**
+ **Linearizable**
+ **Local reads**

**RC protocols:**
- **Relaxed Consistency**
+ **Asynchronous** + **Local reads**

**The L²AW theorem**

Any *Linearizable Asynchronous* read/write register implementation that <u>tolerates a crash</u> (**W**ithout blocking reads or writes), has <u>no **L**ocal reads</u>.

So can we not improve read performance without compromises?

**L²AW vs. CAP**

= Both Linearizability & Asynchrony

🔑 L²AW read performance in its tradeoff
Key for read-dominant workloads

**Fault-tolerance**
- **CAP**: <u>network partitions</u>
  + msg loss + partitioned nodes
  exec ops to violate safety
- **L²AW**: <u>server crashes</u>
  + no msg loss + crashed nodes
  do not exec ops to violate safety

**When must compromise?**
- **CAP**: during network partitions
  (not during partition-free)
  sacrifice safety or progress of ops
- **L²AW**: always sacrifice local reads
  (even if crashes have not occurred)

## Practice

**Almost Local Reads (ALRs)**

Inevitably ALR latency > local reads
💡 But **little or no extra** network and
processing **costs to remote replicas**

**ALRs batch reads with a twist**
🌀 Exec all reads in batch w/ local replica
+ one sync per batch on remote nodes

**Syncs are cheap!**
- writes act as implicit zero-cost syncs
- explicit sync has small constant cost
- 1 sync per batch regardless its size

**Add missing piece to protocols
of all 3 (RC, LS, RA) categories**

| example of reads invoked by a replica | RC | LS | RA | ALRs (this work) |
|---|---|---|---|---|
| read₁(x) | local | local | remote | local |
| read₂(y) | local | local | remote | local |
| ... | ... | ... | ... | local |
| readₙ(z) | local | local | remote | local / sync |
| Linearizable | ✗ | ✓ | ✓ | ✓ |
| Asynchronous | ✓ | ✗ | ✓ | ✓ |
| Cost on remote replicas (network / compute) | zero | zero | O(n) even with traditional batching | small constant- independent of reads in ALR batch → zero when a write is timely |

local: execution uses only local replica
remote: execution involves remote replicas

✓ **RC** with ALRs → **Linearizable**
✓ **LS** with ALRs → **Asynchronous**
✓ **RA** with ALRs → **Performant**

**ALR-enhanced throughput
of state-of-the-art protocols**



+ Asynchronous
+ Linearizable
+ 2x perf

Legend: vanilla protocols / with ALRs
Y-axis: Million Requests per Second (0–250)
X-axis: Hermes (LS), ZAB (RC), Raft (RA)

95% reads | 8B keys 32B vals | 5x R320 Cloudlab nodes (replicas)