Beyond reCAP: Local Reads and Linearizable Asynchronous Replication

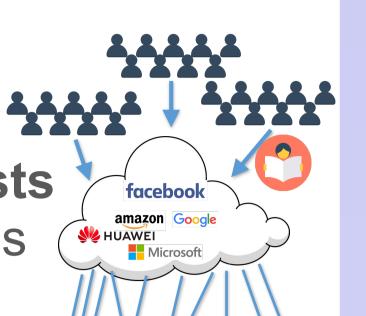
A. Katsarakis*†, E. Gioratmis**, V. Gavrielatos†, P. Bhatotia*, A. Dragojevic⁺, B. Grot*, V. Nagarajan*, P. Fatourou▼ † Huawei Research, *TU Munich, *Citadel Securities, *University of Edinburgh, *University of Crete and FORTH, *Equal contribution

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

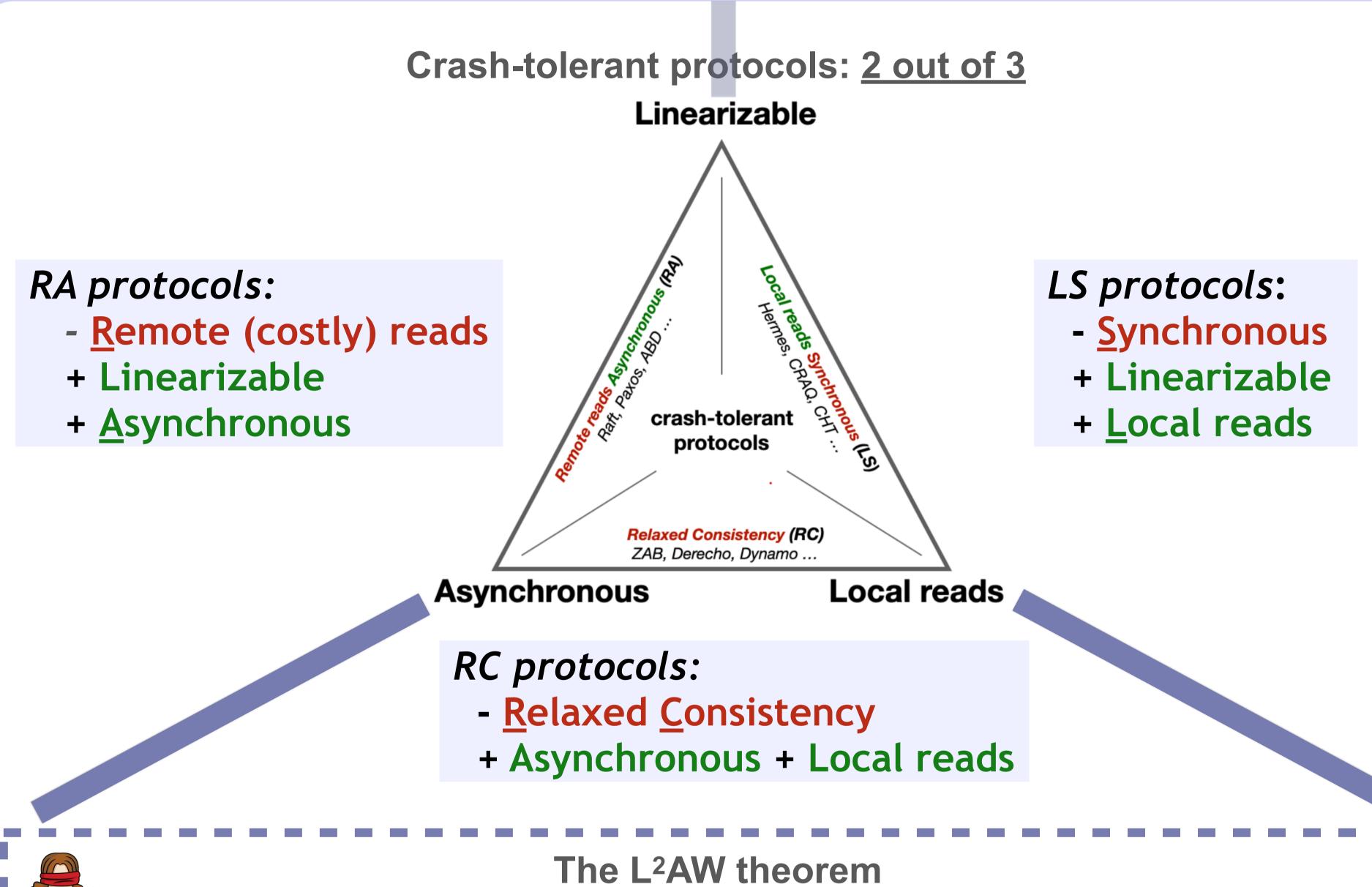


- 1. Linearizable
- 2. Asynchronous
- 3. Local reads: for max perf.





Theory



Any <u>Linearizable Asynchronous</u> read/write register implementation that tolerates a crash (Without blocking reads or writes), has no <u>Local reads</u>. tolerates a crash (Without blocking reads or writes), has no Local reads.



So can we not improve read performance without compromizes?



L²AW vs. CAP



Both Linearizability & Asynchrony

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

Fault-tolerance

CAP: network partitions + msg loss + partitioned nodes exec ops to violate safety

L2AW: server crashes

+ 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)

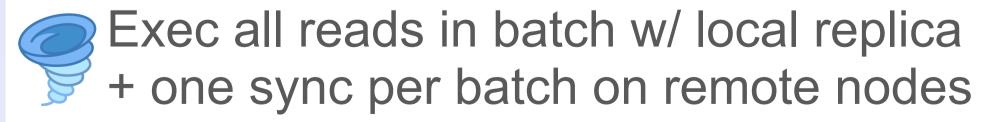


Practice

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



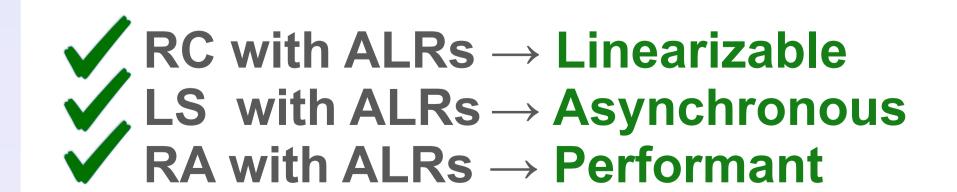
Syncs are cheap!



- writes act as implicit zero-cost syncs
- writes act as implicit zero-cost syncs
 explicit sync has small constant cost
 1 sync per batch regardless its size
 - 1 sync per batch regardless its size



example of reads invoked by a replica	RC	LS	RA	ALRs		
read ₁ (x) read ₂ (y) read _n (z)	local local local		remote remote remote	ALR batch So o o o o o o o o o o o o o o o o o o		eager ocal read (prior/aft
Linearizable	×	/	✓			
Asynchronous	<	×	/			
Cost on remote replicas (network / compute)	zero	zero	O(n)		onstant	
local: execution uses of remote: execution involved	even with traditional batching	independent of reads in ALR batch		when a write is <i>timely</i>		



ALR-enhanced throughput of state-of-the-art protocols vanila protocols 200 with ALRs + Linearizable + 2x perf Hermes (LS) ZAB (RC) Raft (RA)

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



†This work occurred when the authors were at the University of Edinburgh











