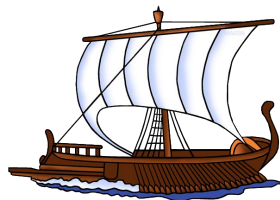




Odyssey: The Impact of Modern Hardware on Strongly-Consistent Replication Protocols

Vasilis Gavrielatos, Antonios Katsarakis, Vijay Nagarajan



Thanks to:



Scope

- Replication protocols
- Strong consistency
- Replicated KVS
- Read/Write API
- Datacenter

Scope

- Replication protocols
- Strong consistency
- Replicated KVS
- Read/Write API
- Datacenter



Paxos
Raft
Zookeeper
Chain replication

Scope

- Replication protocols
- Strong consistency
- Replicated KVS
- Read/Write API
- Datacenter

Old Hardware

- Single-thread
- Slow Disk
- Slow network

Scope

- Replication protocols
- Strong consistency
- Replicated KVS
- Read/Write API
- Datacenter

~~Old~~ Modern Hardware

- ~~Single thread~~ Manycore
- ~~Slow Disk~~ Big Memories
- ~~Slow network~~ Fast Networks

Scope

- Replication protocols
- Strong consistency
- Replicated KVS
- Read/Write API
- Datacenter

~~Old~~ Modern Hardware

- ~~Single thread~~ Manycore
- ~~Slow Disk~~ Big Memories
- ~~Slow network~~ Fast Networks

Scope

- Replication protocols
- Strong consistency
- Replicated KVS
- Read/Write API
- Datacenter

~~Old~~ Modern Hardware

- ~~Single thread~~ Manycore
- ~~Slow Disk~~ Big Memories
- ~~Slow network~~ Fast Networks

Scope

- Replication protocols
- Strong consistency
- Replicated KVS
- Read/Write API
- Datacenter

~~Old~~ Modern Hardware

- ~~Single thread~~ Manycore
- ~~Slow Disk~~ Big Memories
- ~~Slow network~~ Fast Networks

Scope

- Replication protocols
- Strong consistency
- Replicated KVS
- Read/Write API
- Datacenter

~~Old~~ Modern Hardware

- ~~Single thread~~ Manycore
- ~~Slow Disk~~ Big Memories
- ~~Slow network~~ Fast Networks

Modern Hardware challenges conventional wisdom

How do protocols perform over modern hardware?

How do protocols perform over modern hardware?

What are the best practices?

Taxonomy

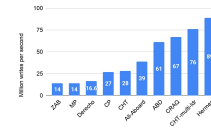
Leader-based Total Order	Leader-based Per-key Order
Decentralized Total Order	Decentralized Per-key Order



Odyssey Framework



Design Space Characterization



Taxonomy

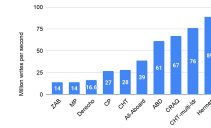
Leader-based Total Order	Leader-based Per-key Order
Decentralized Total Order	Decentralized Per-key Order



Odyssey Framework



Design Space Characterization



Taxonomy

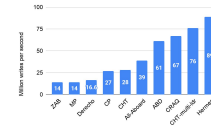
Leader-based Total Order	Leader-based Per-key Order
Decentralized Total Order	Decentralized Per-key Order



Odyssey Framework



Design Space Characterization



Taxonomy

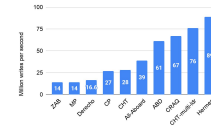
Leader-based Total Order	Leader-based Per-key Order
Decentralized Total Order	Decentralized Per-key Order



Odyssey Framework



Design Space Characterization

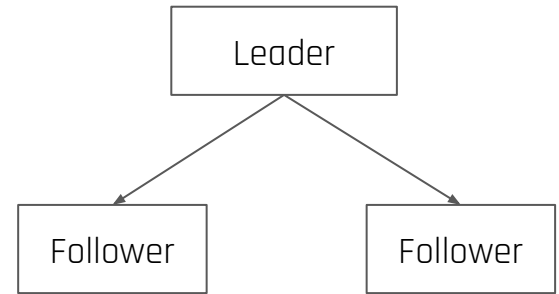
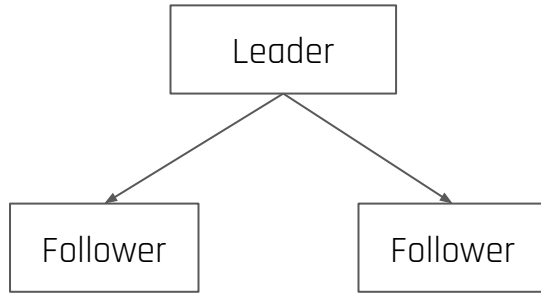


	Total Order	Per-key Order
Leader-based		
Decentralized		

Total Order

Per-key Order

Leader-based

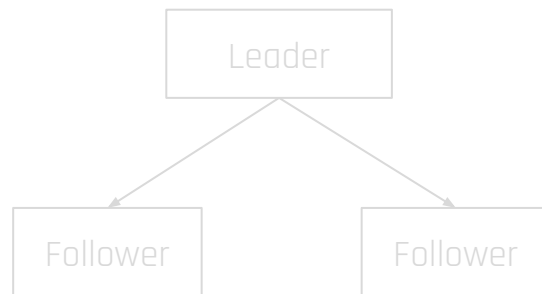
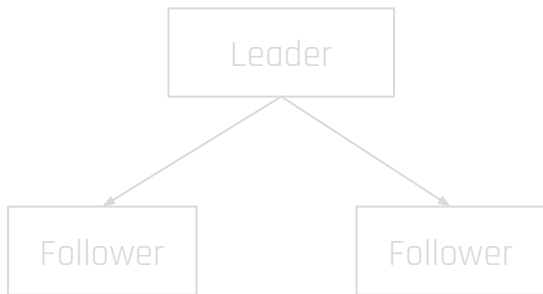


Decentralized

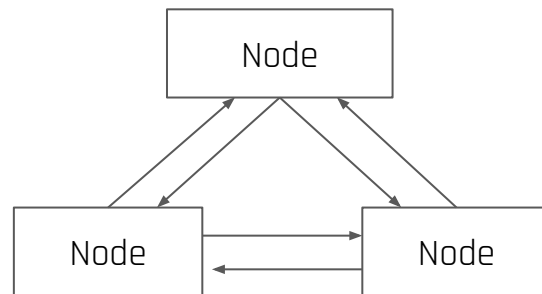
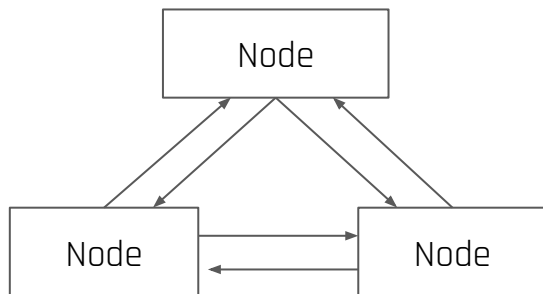
Total Order

Per-key Order

Leader-based



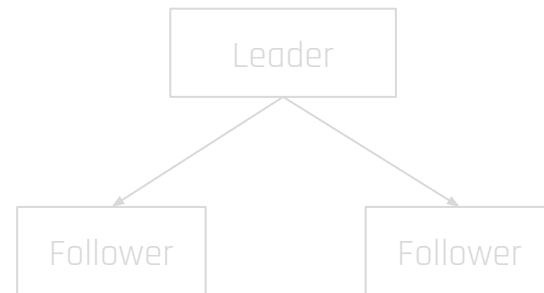
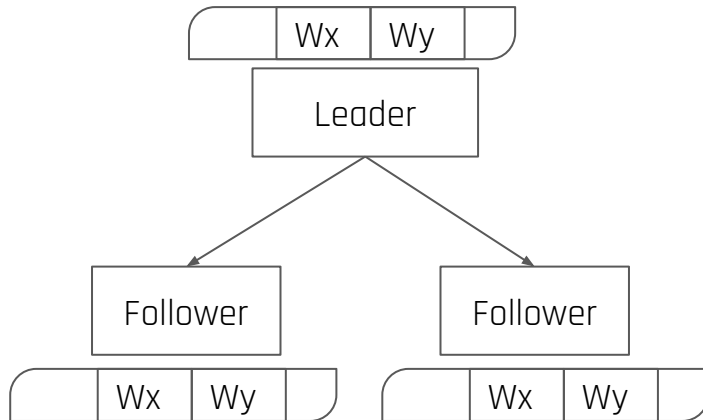
Decentralized



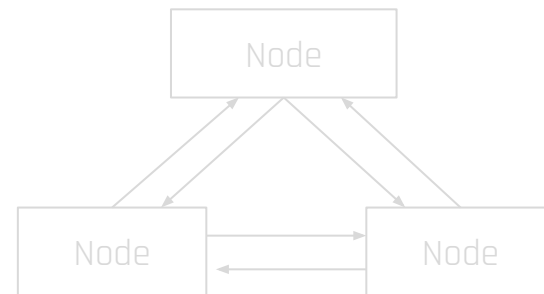
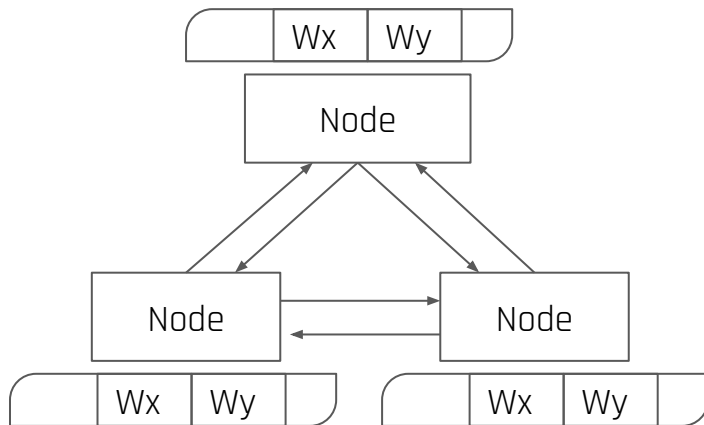
Total Order

Per-key Order

Leader-based



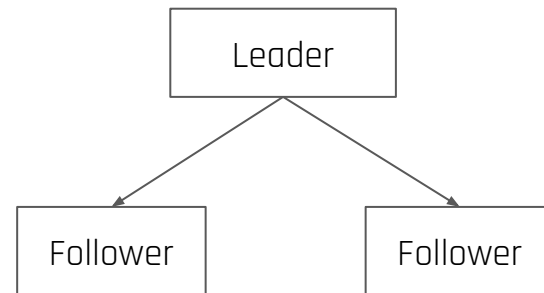
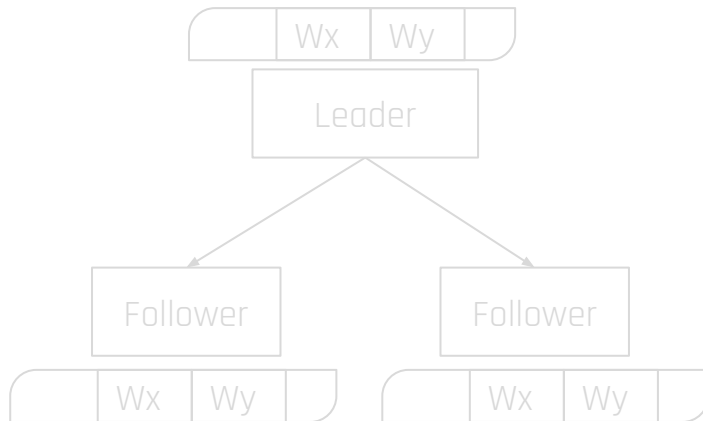
Decentralized



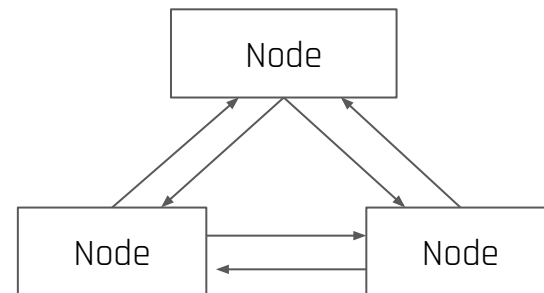
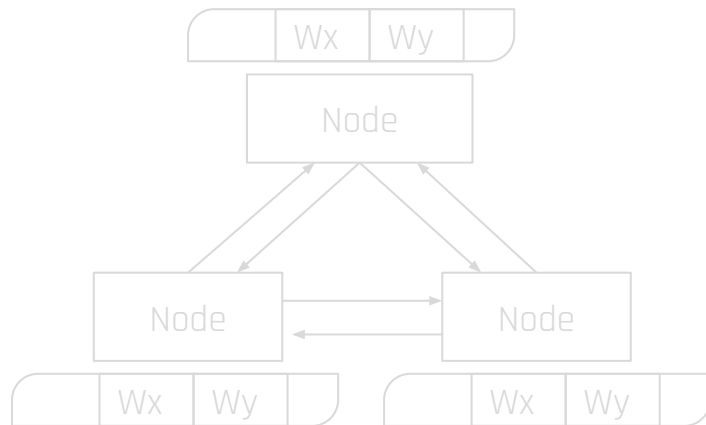
Total Order

Per-key Order

Leader-based



Decentralized



	Total Order	Per-key Order
Leader-based	<p>Multi-Paxos (MP)</p> <p>↓</p> <p>ZAB</p>	<p>CHT</p> <p>↙ ↘</p> <p>CRAQ multi-ldr CHT</p>
Decentralized	<p>Derecho</p>	<p>Classic Paxos (CP)</p> <p>↙ ↓ ↘</p> <p>ABD Hermes All-aboard Paxos</p>

Taxonomy

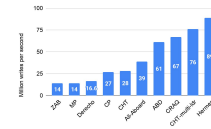
Leader-based Total Order	Leader-based Per-key Order
Decentralized Total Order	Decentralized Per-key Order



Odyssey Framework



Design Space Characterization



Servers	5
Hardware threads	40
Network Bandwidth	56 Gbps (RDMA)
KVS	MICA
Key-value size	48 B

General Directives.

- Prioritize thread-scalability, then load-balance and then the work-per-request ratio. Total order should be avoided in read/write systems.
- Leader-based protocols can achieve high-performance, but care must be taken to ensure load balance.
- It is worth investing in the hardware multicast primitive only in the case of LPK0 protocols.
- Local reads can deliver great performance, but it's not guaranteed.
- In order to minimize latency, choose protocols with high throughput.

Recommendations

- All-aboard is the most attractive design point for a scenario where: 1) availability is the most important concern and 2) conditional writes are required.
- If simple writes will do, then we recommend ABD.
- If a small window of unavailability on a failure is tolerable, then Hermes is the best candidate, while CHT-multi-ldr and CRAQ are good alternatives.

General Directives.

- Prioritize thread-scalability, then load-balance and then the work-per-request ratio. Total order should be avoided in read/write systems.
- Leader-based protocols can achieve high-performance, but care must be taken to ensure load balance.

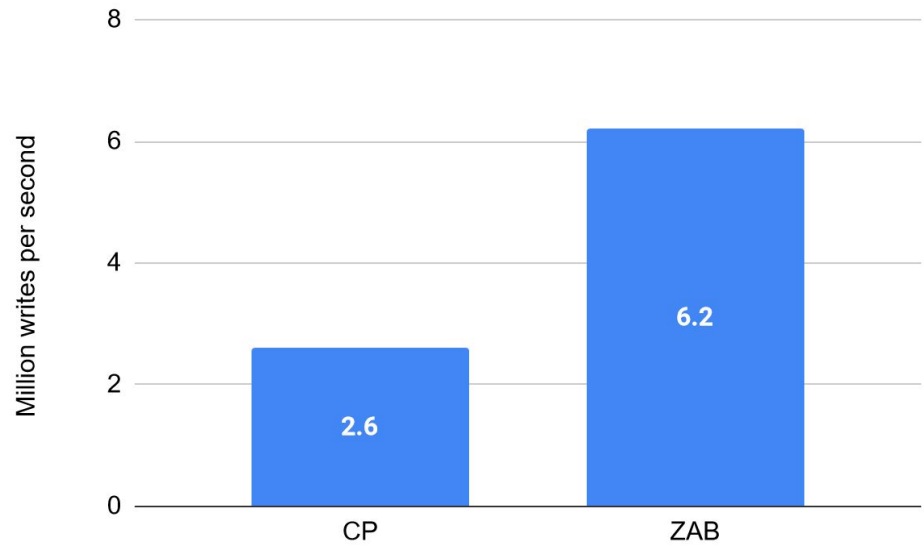
Total order should be avoided in read/write systems.

RECOMMENDATIONS

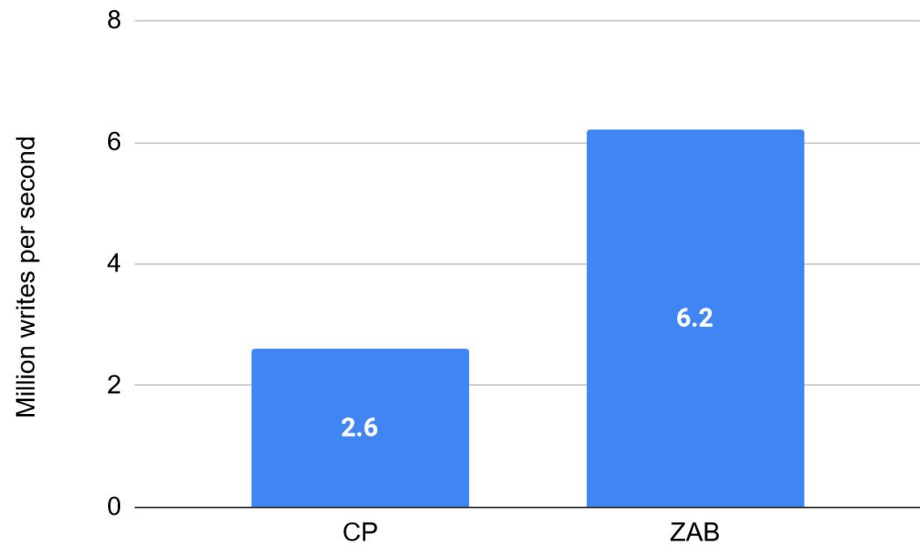
- All-aboard is the most attractive design point for a scenario where: 1) availability is the most important concern and 2) conditional writes are required.
- If simple writes will do, then we recommend ABD.
- If a small window of unavailability on a failure is tolerable, then Hermes is the best candidate, while CHT-multi-ldr and CRAQ are good alternatives.

	Total Order	Per-key Order
Leader-based	<p>Multi-Paxos (MP)</p> <p>↓</p> <p>ZAB</p> <p>[Hunt et al. ATC'10]</p>	<p>CHT</p> <p>↙ ↘</p> <p>CRAQ multi-ldr CHT</p>
Decentralized	<p>Derecho</p>	<p>Classic Paxos (CP)</p> <p>↙ ↓ ↘</p> <p>ABD Hermes All-aboard Paxos</p>

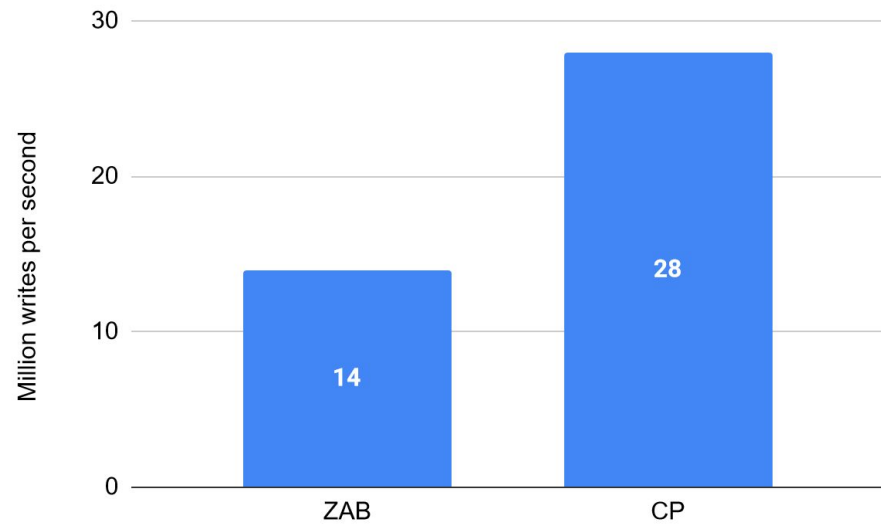
Single-Threaded



Single-Threaded



Multi-Threaded



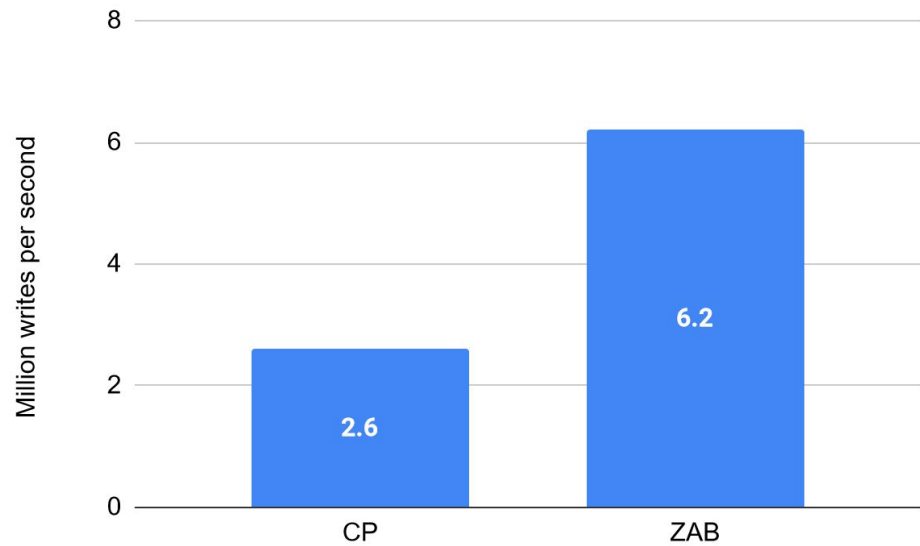
Extending Classic Paxos for High-performance Read-Modify-Write Registers

--All-aboard [Howard's thesis 2019]

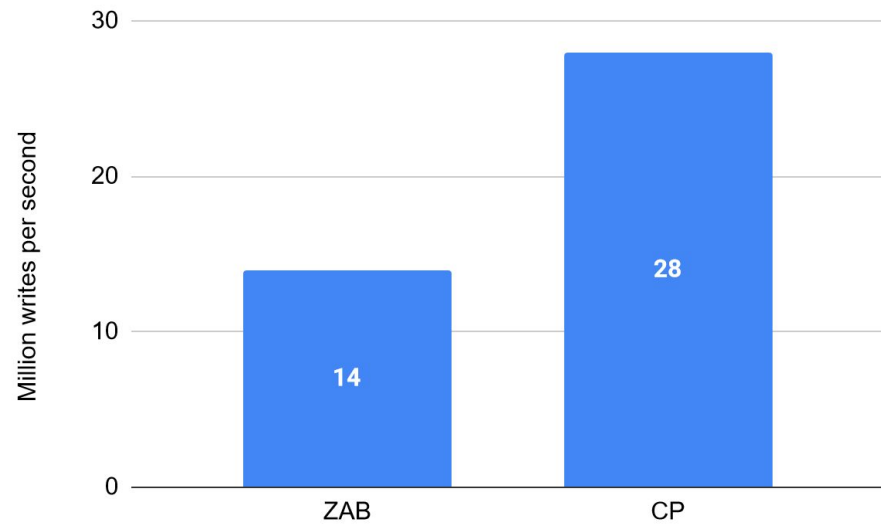
--ABD reads/writes -- Carstamps [Burke NSDI '20]

<https://arxiv.org/abs/2103.14701>

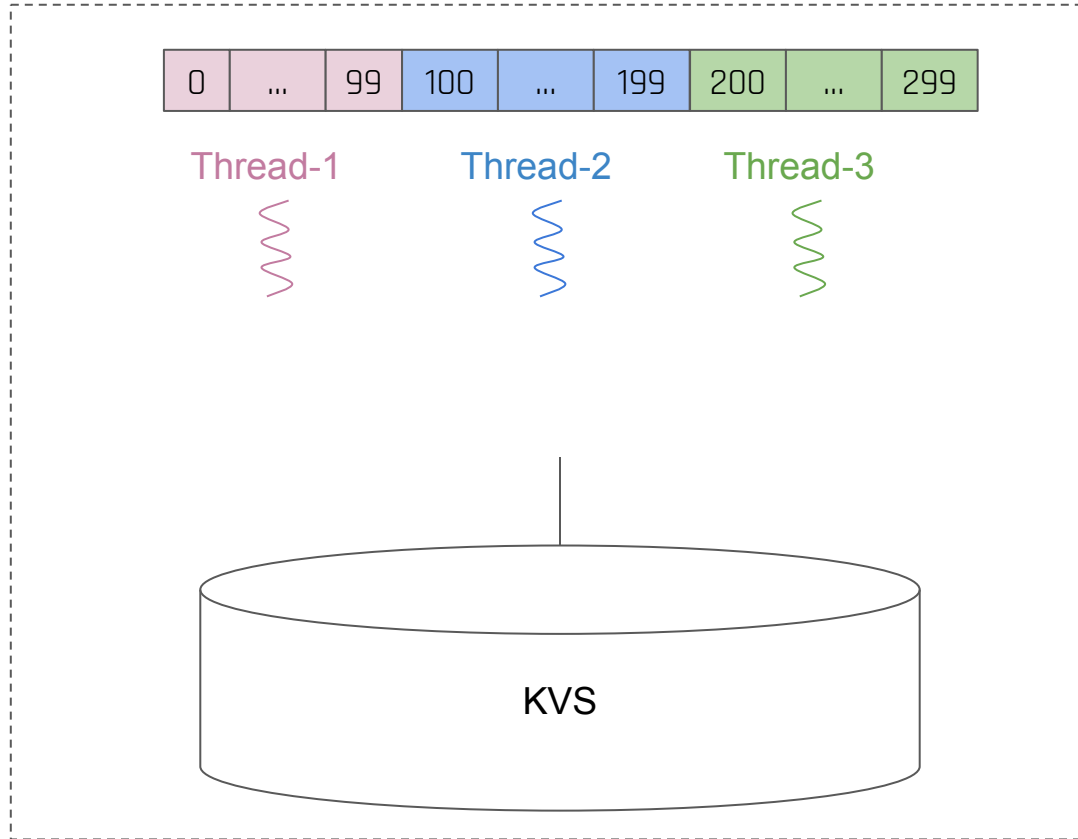
Single-Threaded



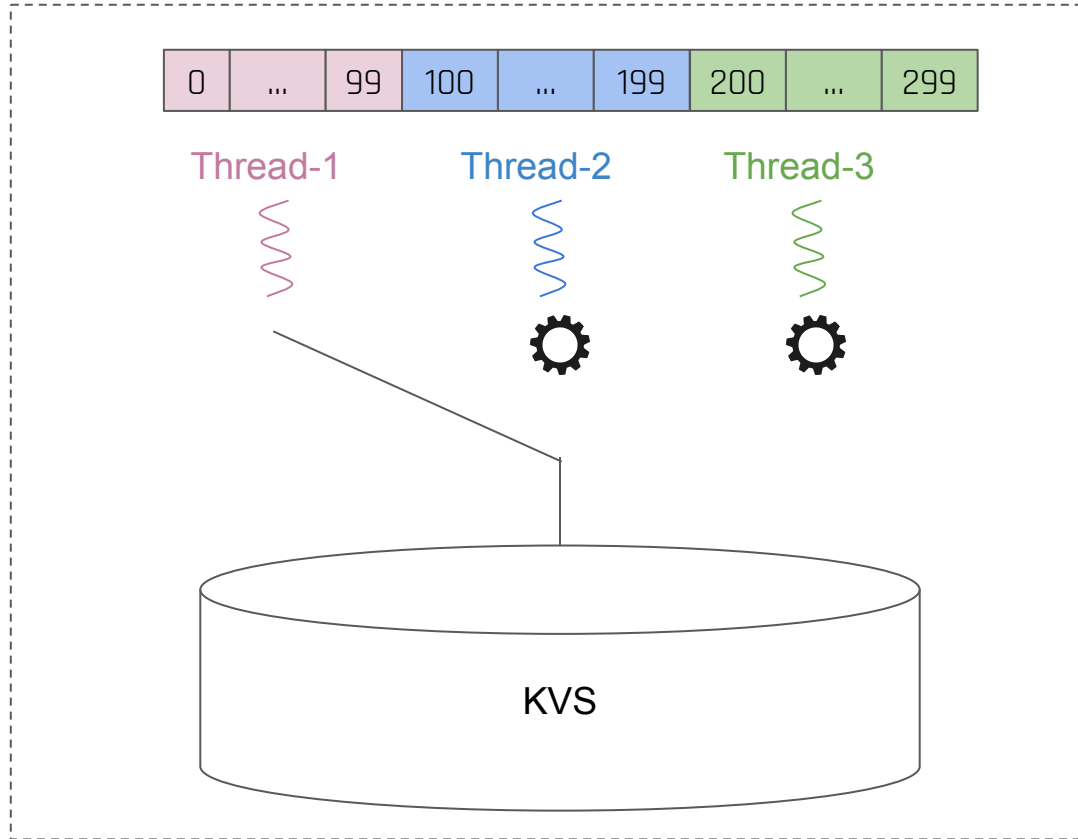
Multi-Threaded



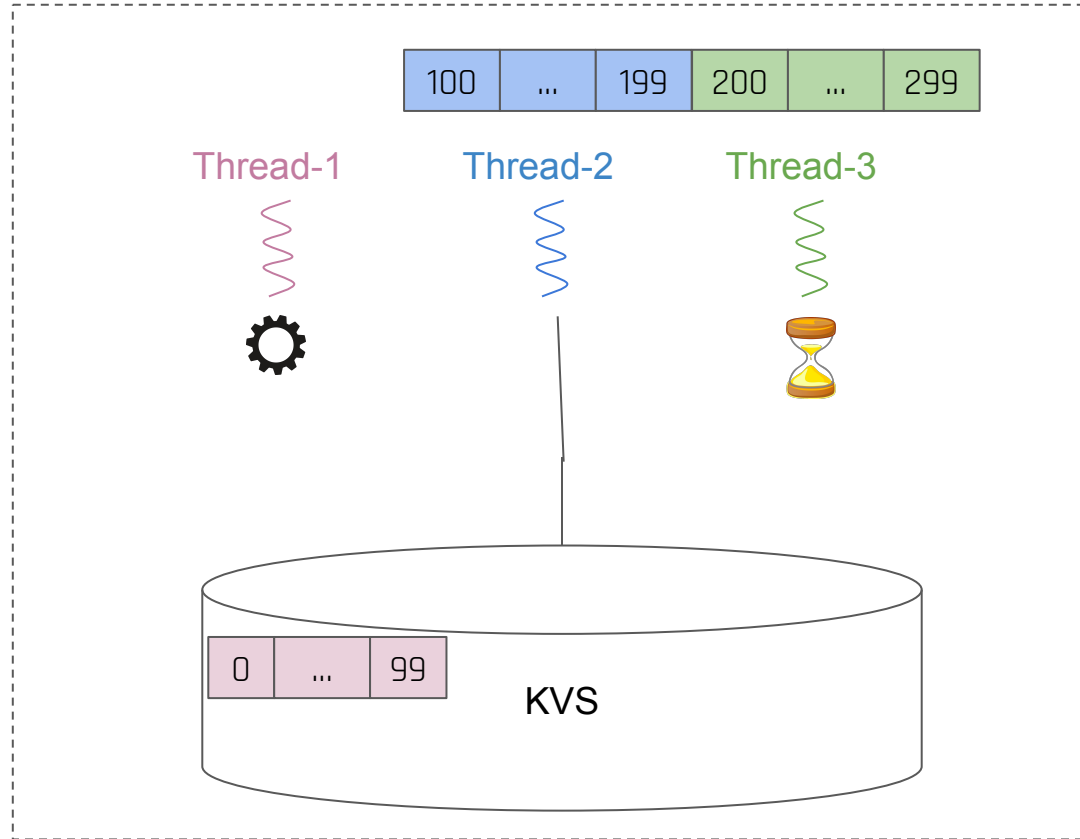
ZAB server



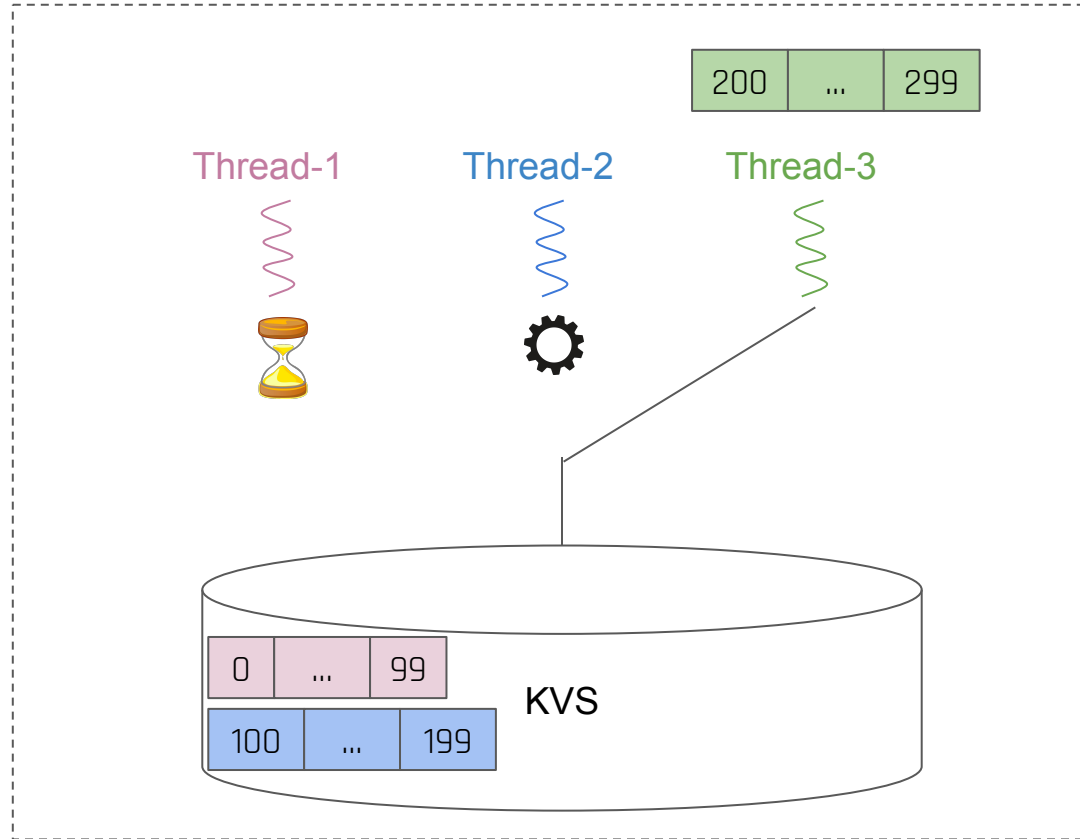
ZAB server



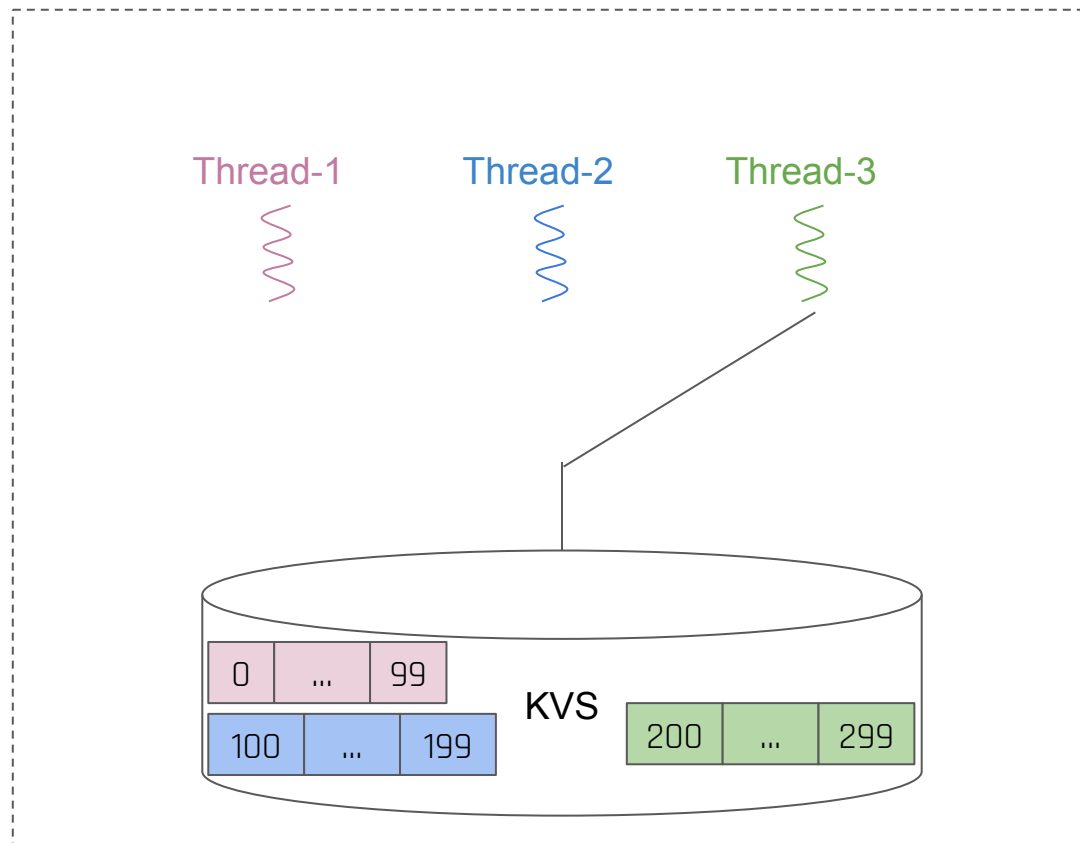
ZAB server

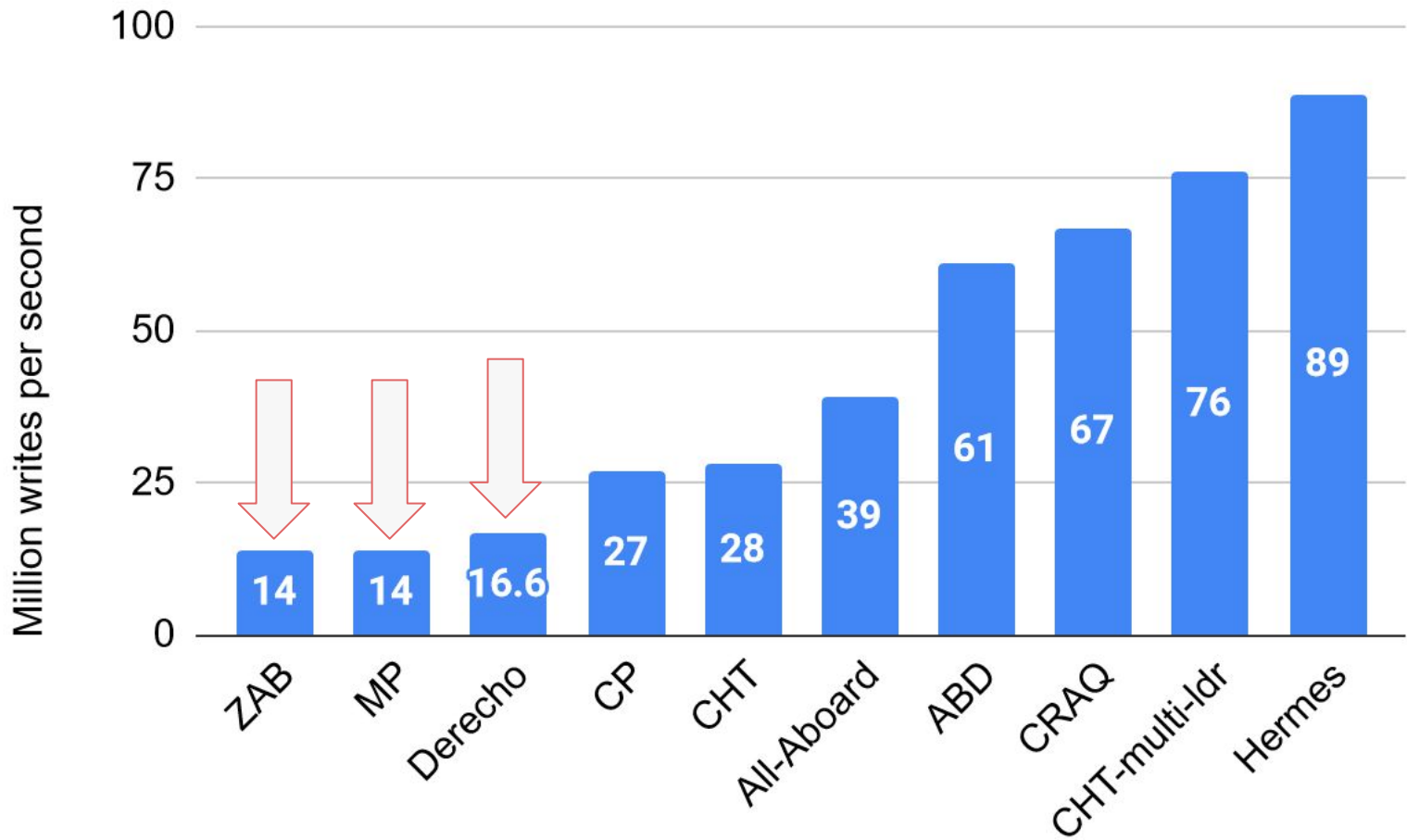


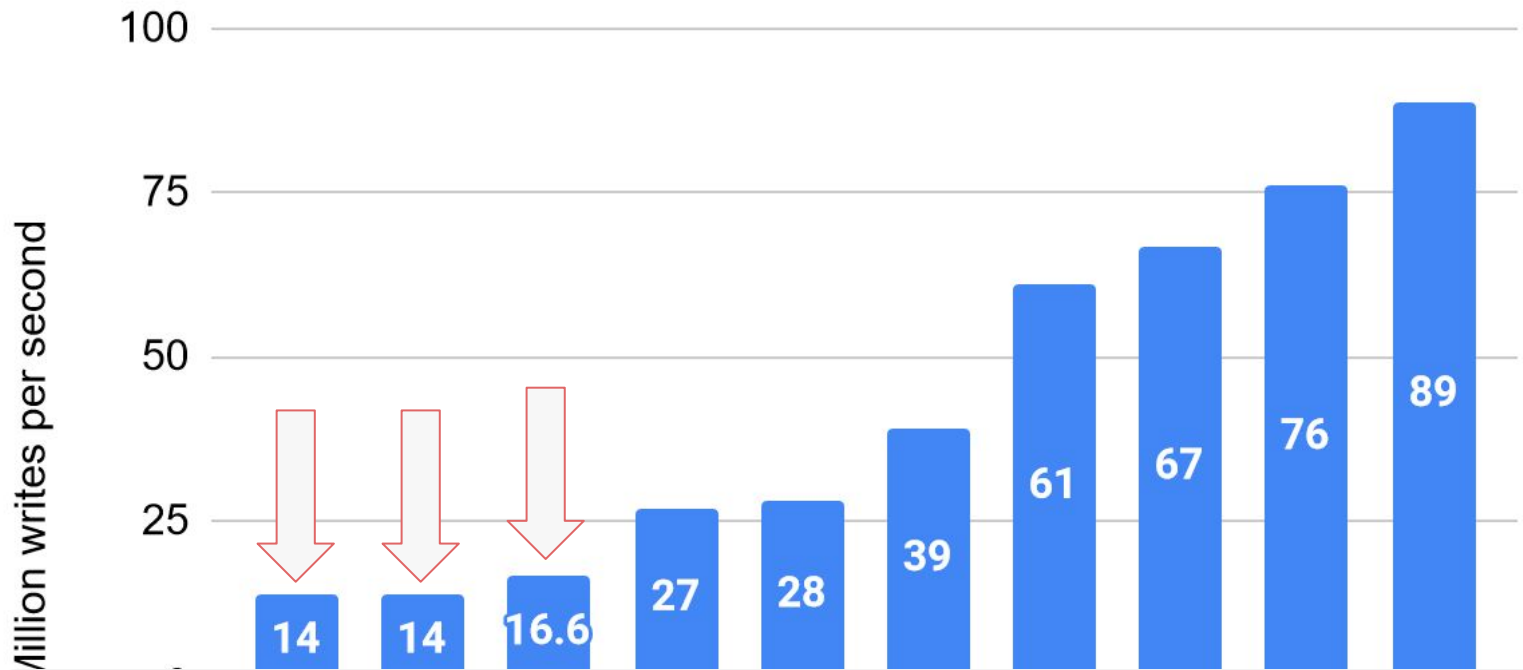
ZAB server



ZAB server







Total-order should be avoided

General Directives.

- Prioritize thread-scalability, then load-balance and then the work-per-request ratio. Total order should be avoided in read/write systems.
- Leader-based protocols can achieve high-performance, but care must be taken to ensure load balance.

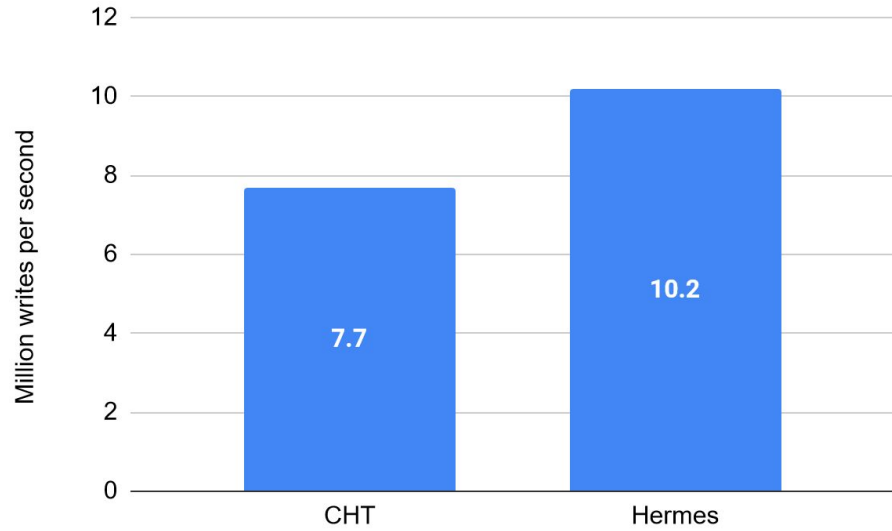
Leader-based protocols can achieve high-performance

RECOMMENDATIONS

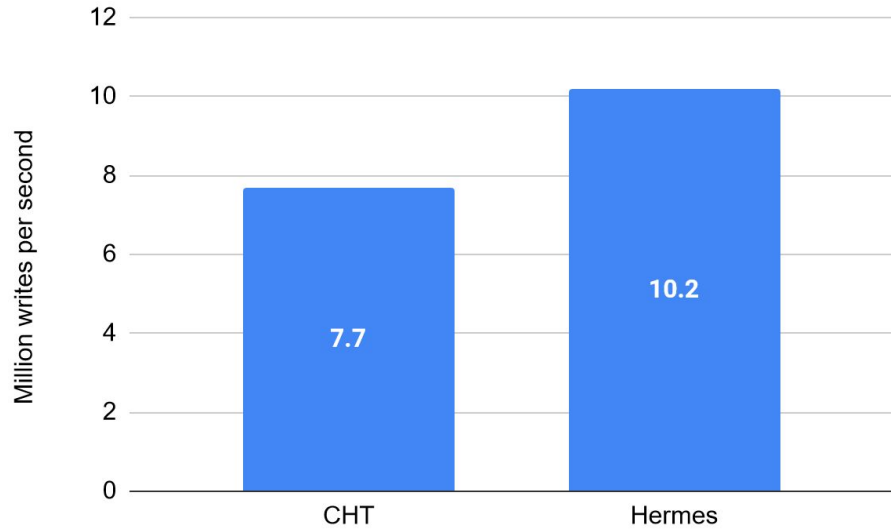
- All-aboard is the most attractive design point for a scenario where: 1) availability is the most important concern and 2) conditional writes are required.
- If simple writes will do, then we recommend ABD.
- If a small window of unavailability on a failure is tolerable, then Hermes is the best candidate, while CHT-multi-ldr and CRAQ are good alternatives.

	Total Order	Per-key Order
Leader-based	<p>Multi-Paxos (MP)</p> <p>↓</p> <p>ZAB</p>	<p>CHT</p> <p>[Chandra et al. PODC '16]</p> <p>↙ ↘</p> <p>CRAQ multi-ldr CHT</p>
Decentralized	<p>Derecho</p>	<p>Classic Paxos (CP)</p> <p>↙ ↓ ↘</p> <p>ABD Hermes All-aboard Paxos</p>

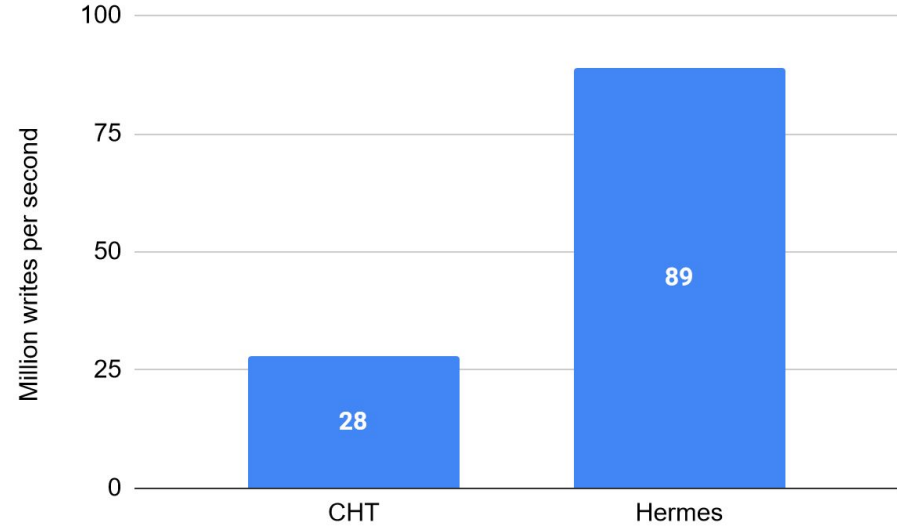
Single-threaded



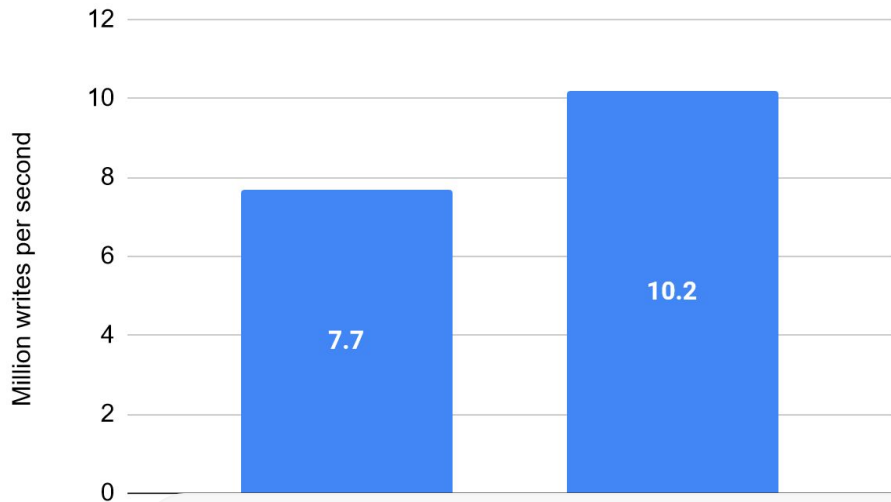
Single-threaded



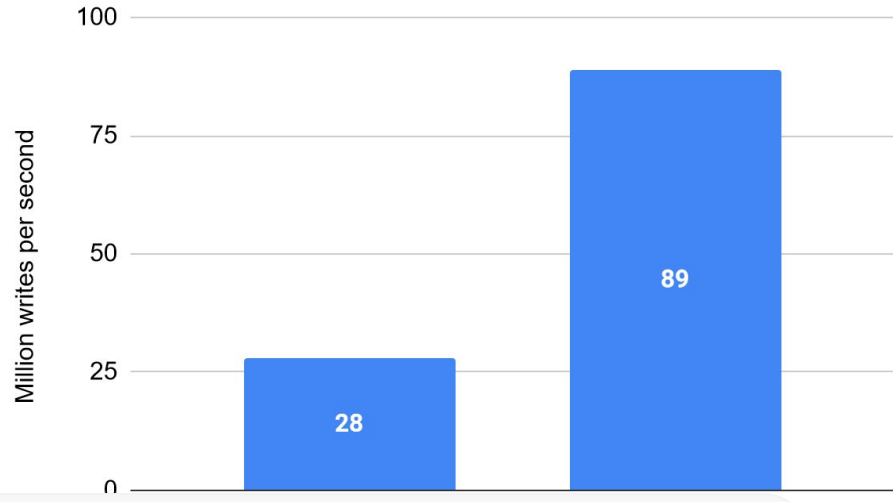
Multi-threaded



Single-threaded

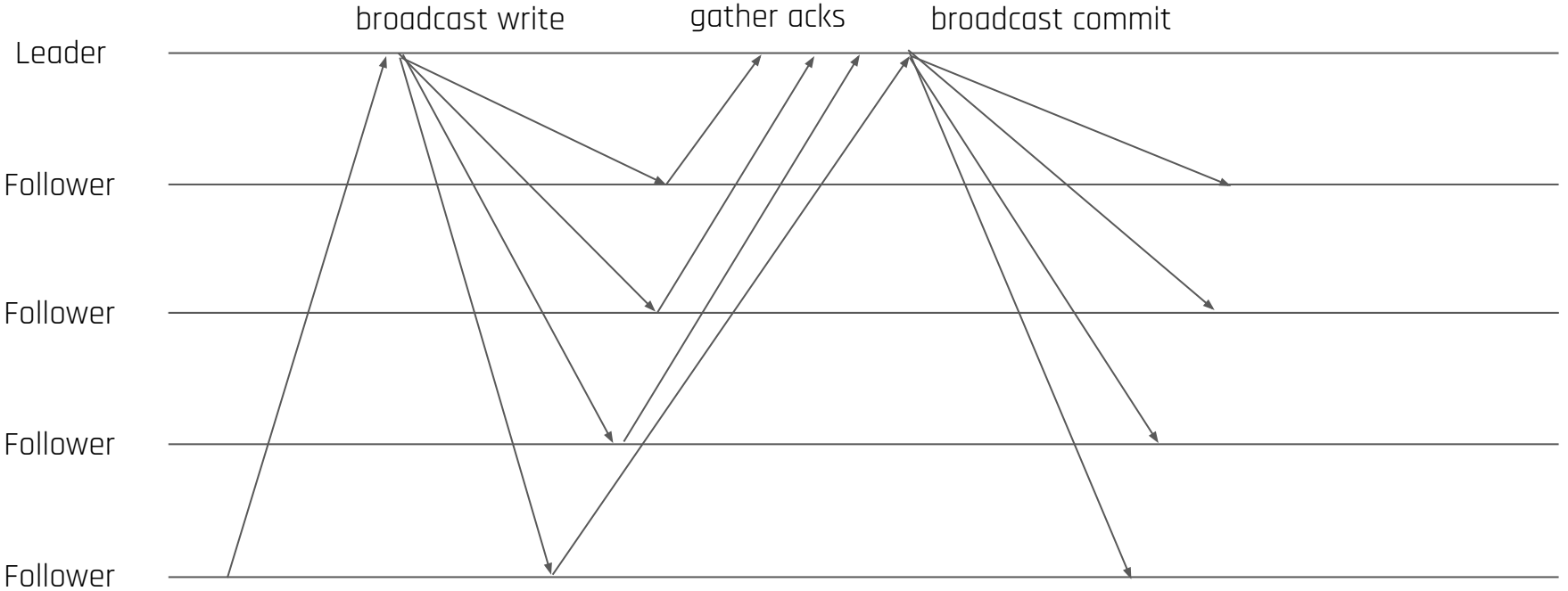


Multi-threaded

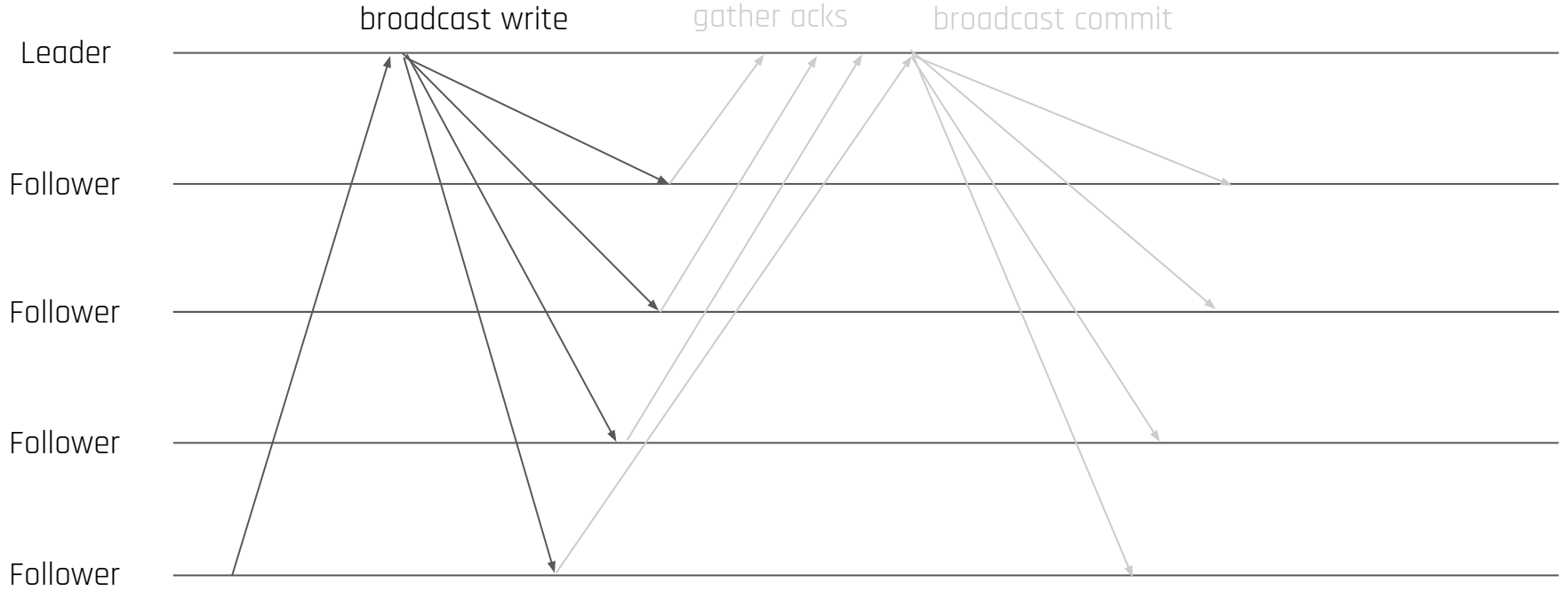


Bottleneck: Leader Send Bandwidth

CHT



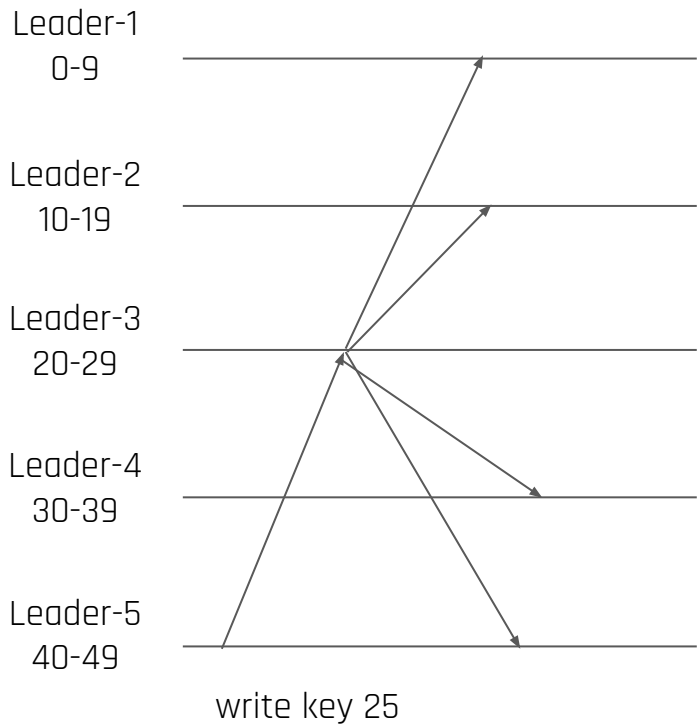
CHT



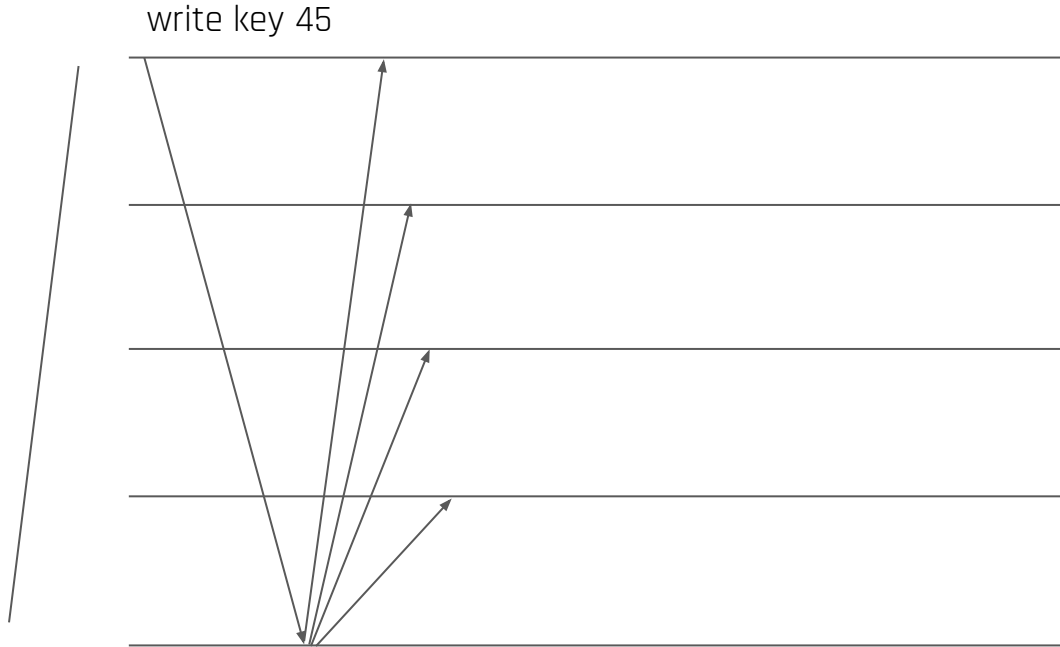
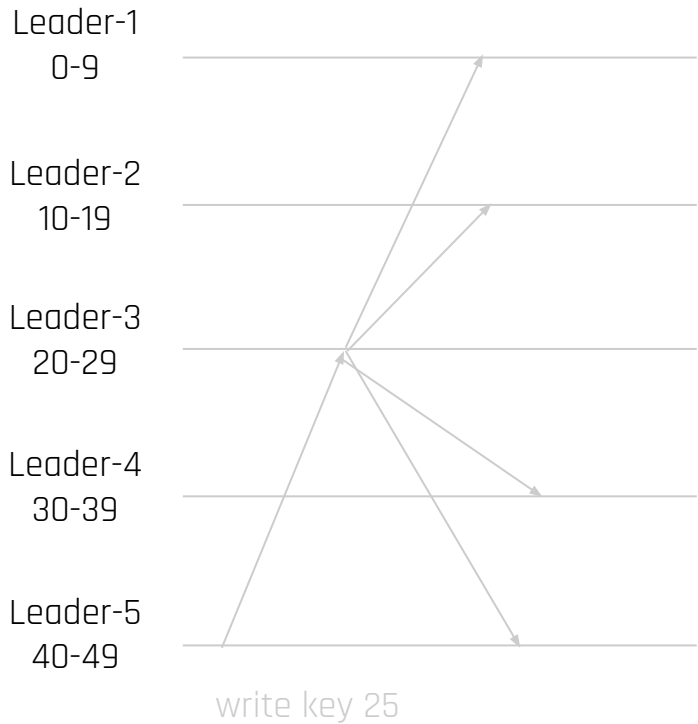
Two Protocol-level Solutions

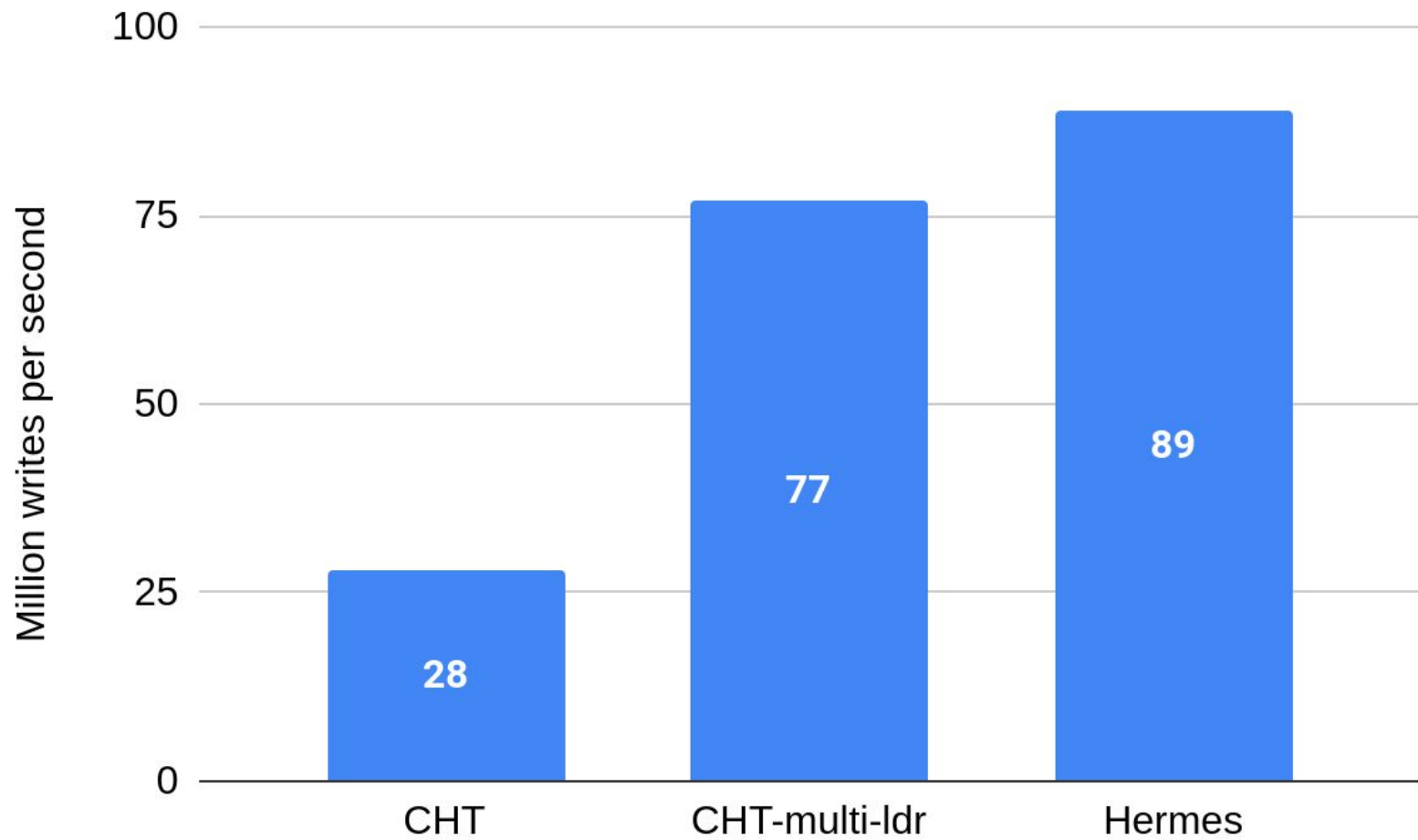
1. Multiple Leaders
2. Chain instead of broadcast

1. Multiple Leaders

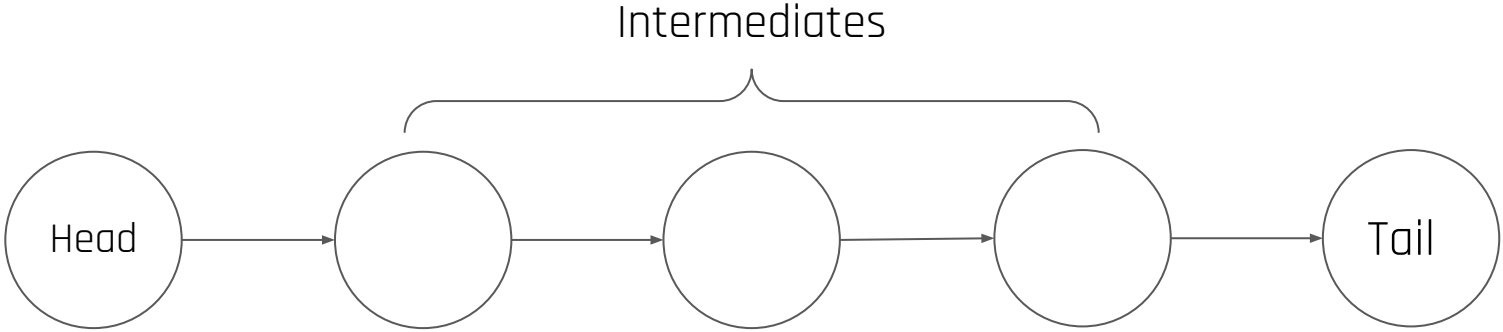


1. Multiple Leaders

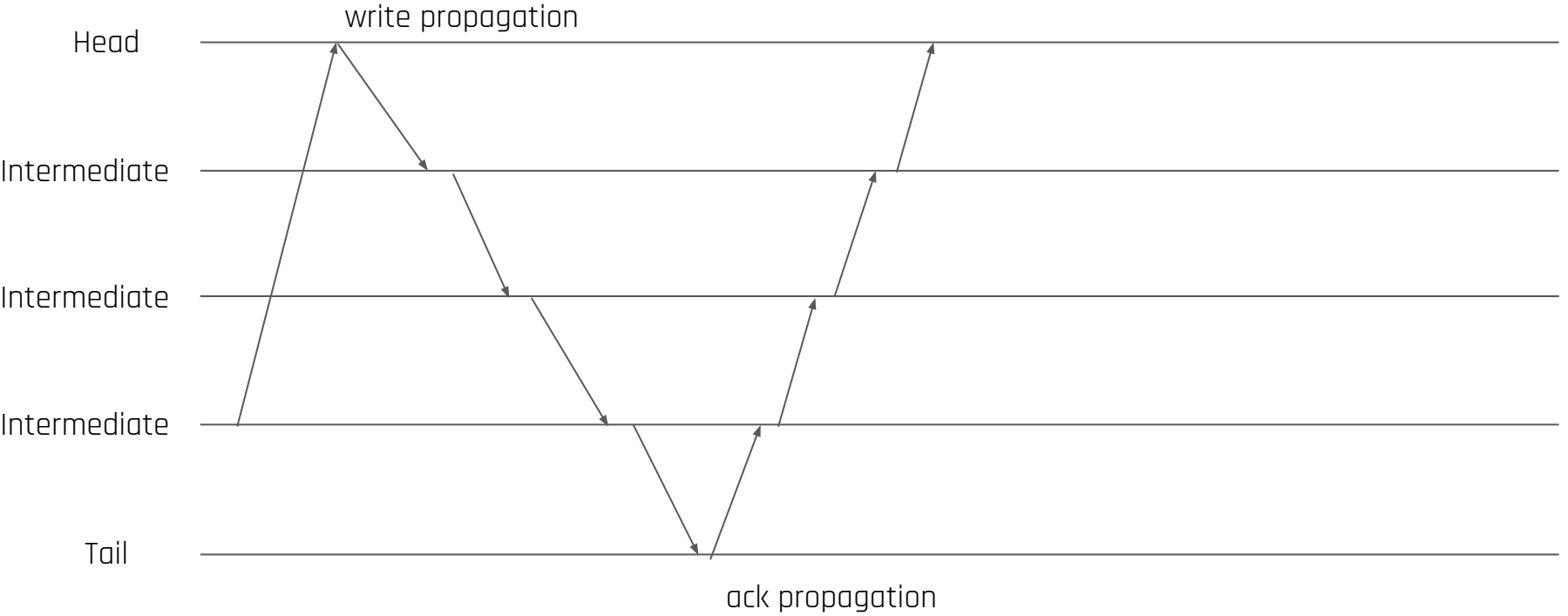


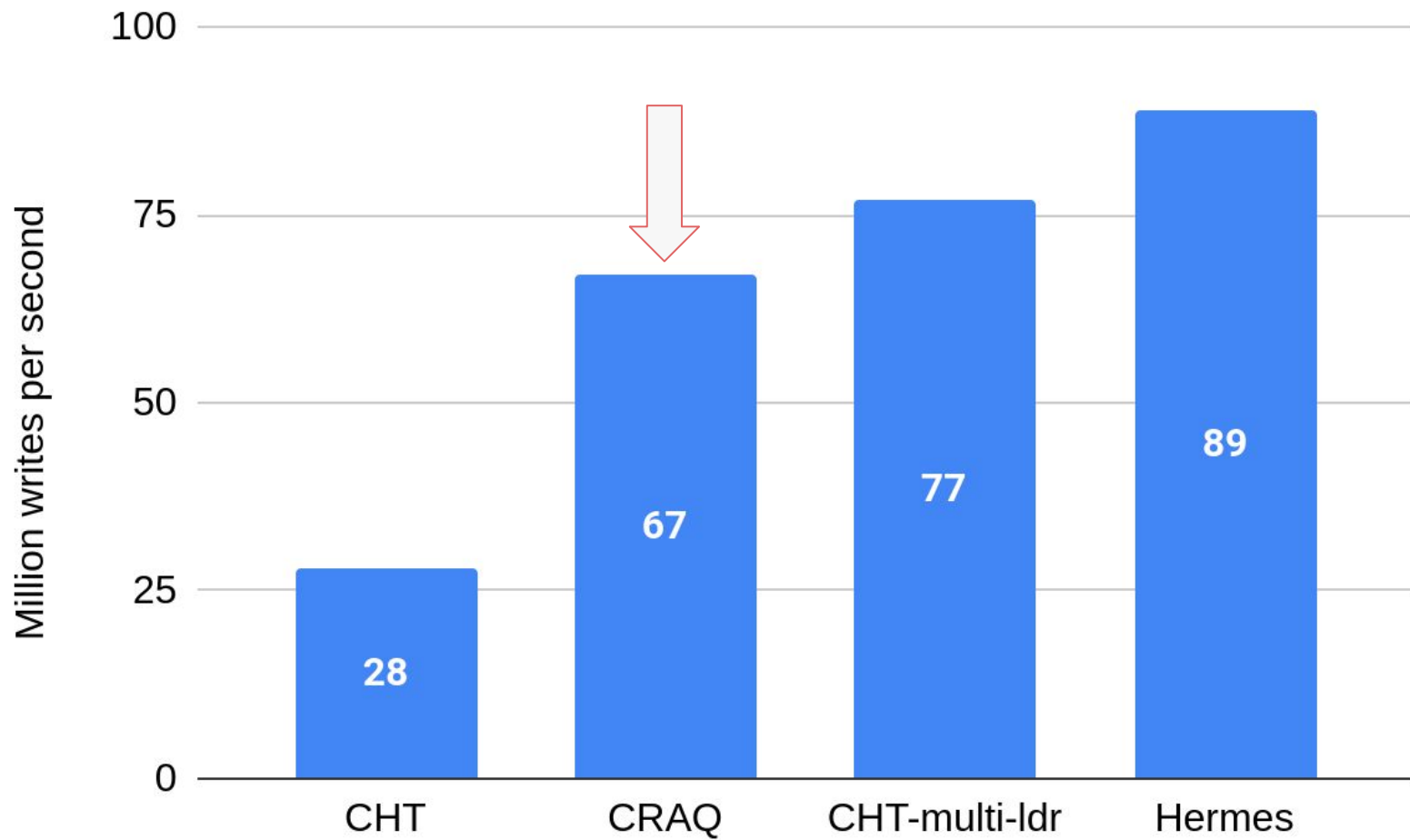


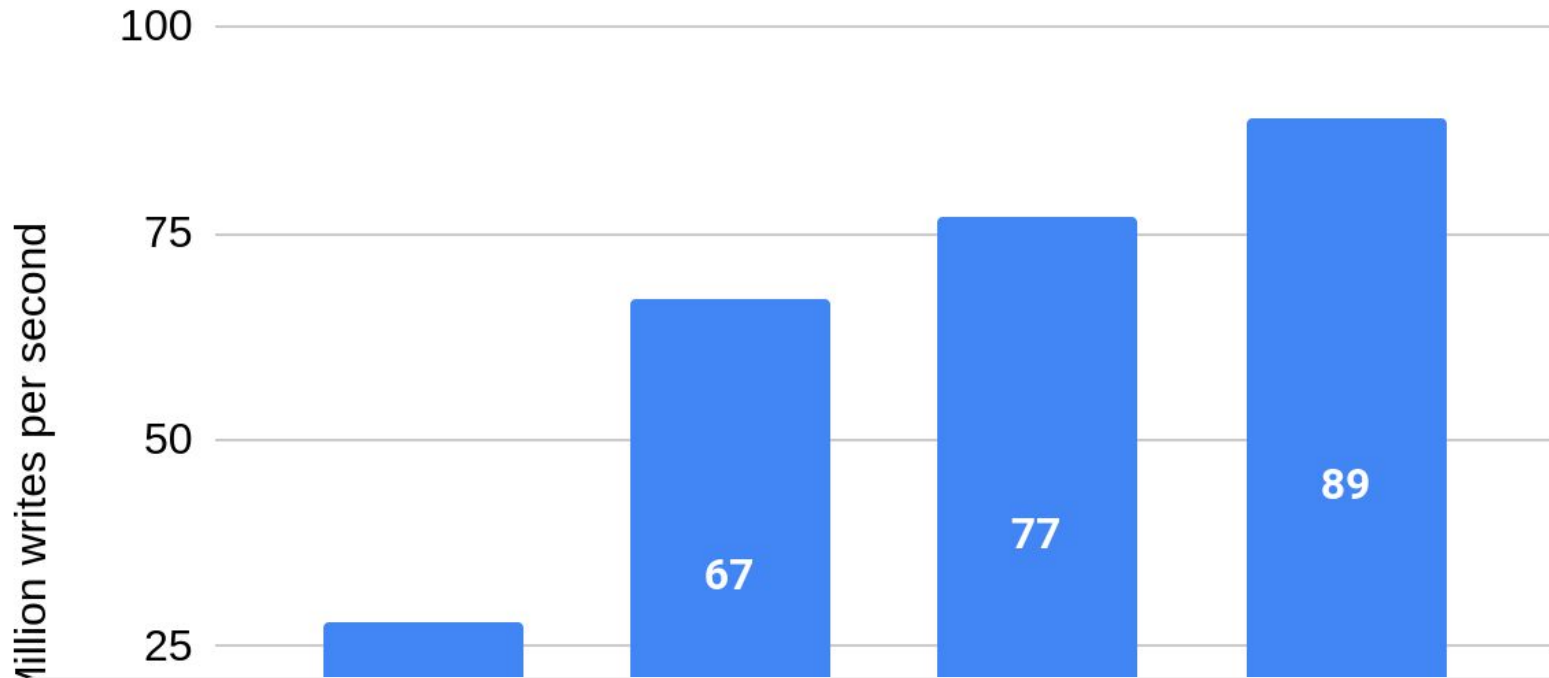
2. Chains instead of Broadcast



2. Chains instead of Broadcast

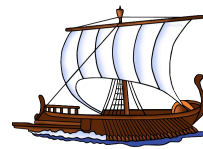






Leader-based can achieve High-Performance

Odyssey Summary



Odyssey Summary



1. Taxonomy

Leader-based Total Order	Leader-based Per-key Order
Decentralized Total Order	Decentralized Per-key Order

- Protocol selection

Odyssey Summary



1. Taxonomy

Leader-based Total Order	Leader-based Per-key Order
Decentralized Total Order	Decentralized Per-key Order

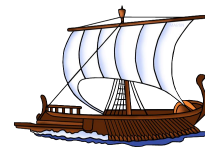
- Protocol selection

2. Odyssey framework



- Fast development of protocols
- Apples-to-apples comparison

Odyssey Summary



1. Taxonomy

Leader-based Total Order	Leader-based Per-key Order
Decentralized Total Order	Decentralized Per-key Order

- Protocol selection

2. Odyssey framework



- Fast development of protocols
- Apples-to-apples comparison

3. Design space Characterization

- Recommendations
- Best practices

