



ShiftLock: Mitigate One-sided RDMA Lock Contention via Handover

Jian Gao, Qing Wang, Jiwu Shu*

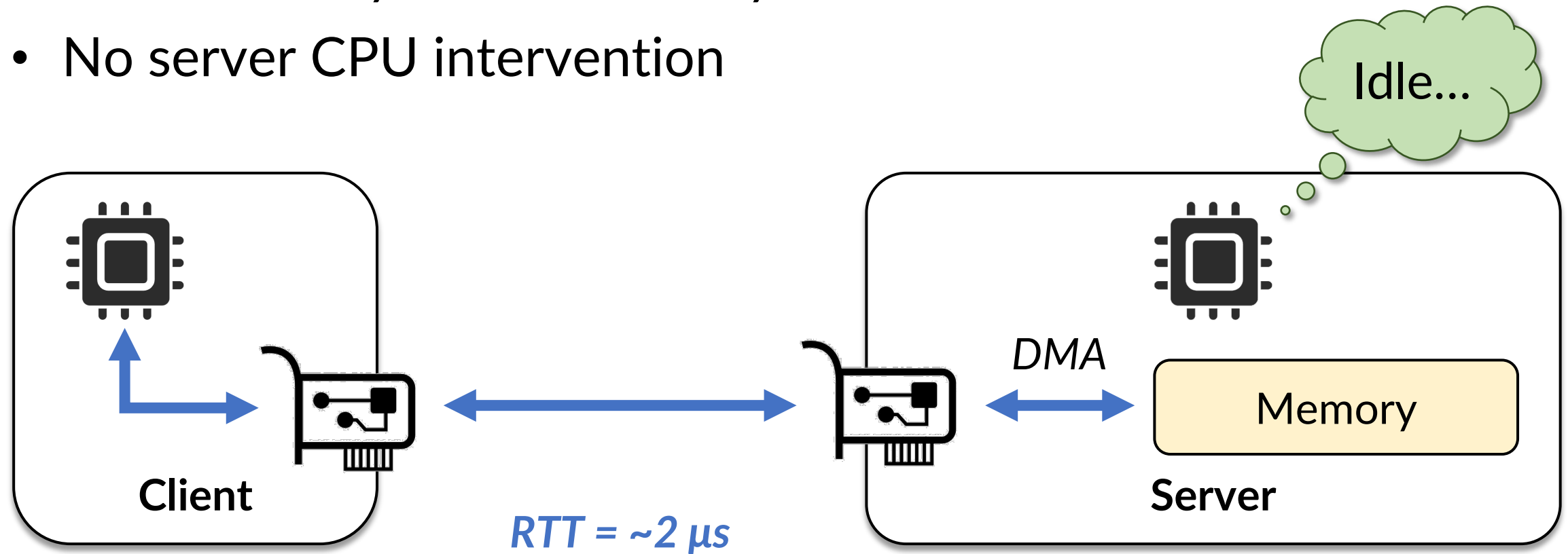
Tsinghua University

2025/02/26 at FAST'25



Remote Direct Memory Access (RDMA)

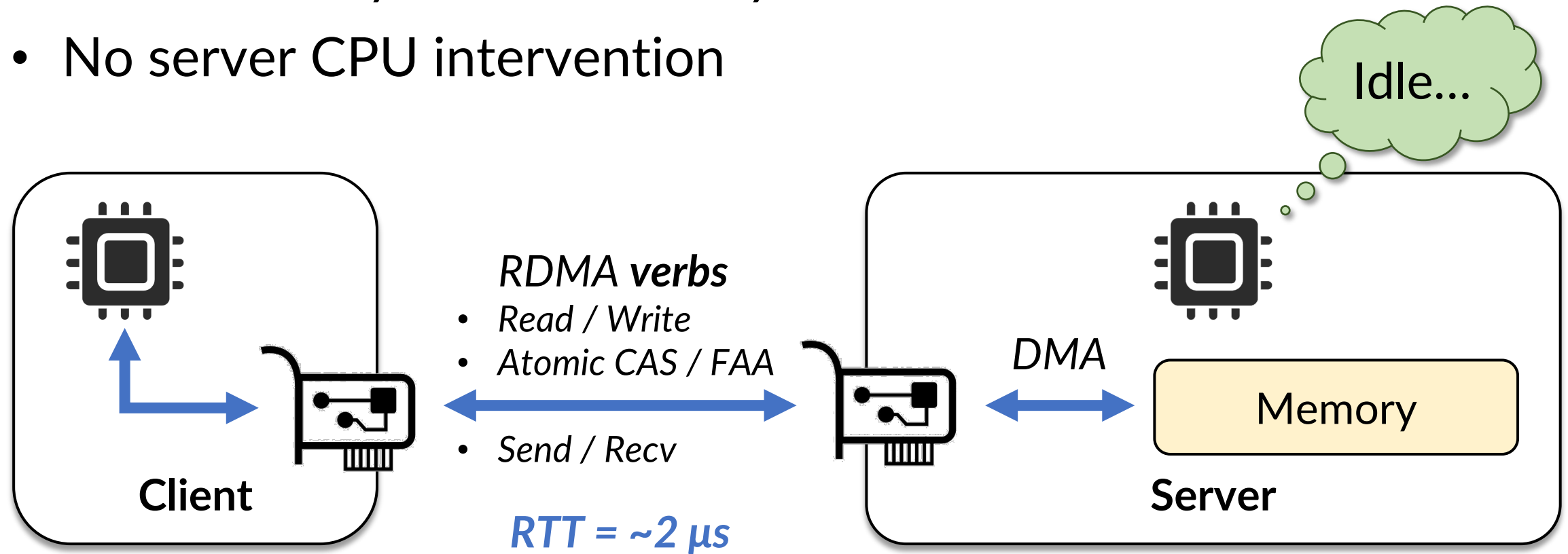
- Client directly access memory on server
- No server CPU intervention





Remote Direct Memory Access (RDMA)

- Client directly access memory on server
- No server CPU intervention





RDMA is Very Powerful

(Extended) Atomic verbs can ...

- Operate on larger-than-8B values
 - *16-byte atomics* with similar performance to original 8-byte atomics
- Compare-swap (CAS) on arbitrary bits
- Fetch-add (FAA) on separate fields



RDMA is Very Powerful

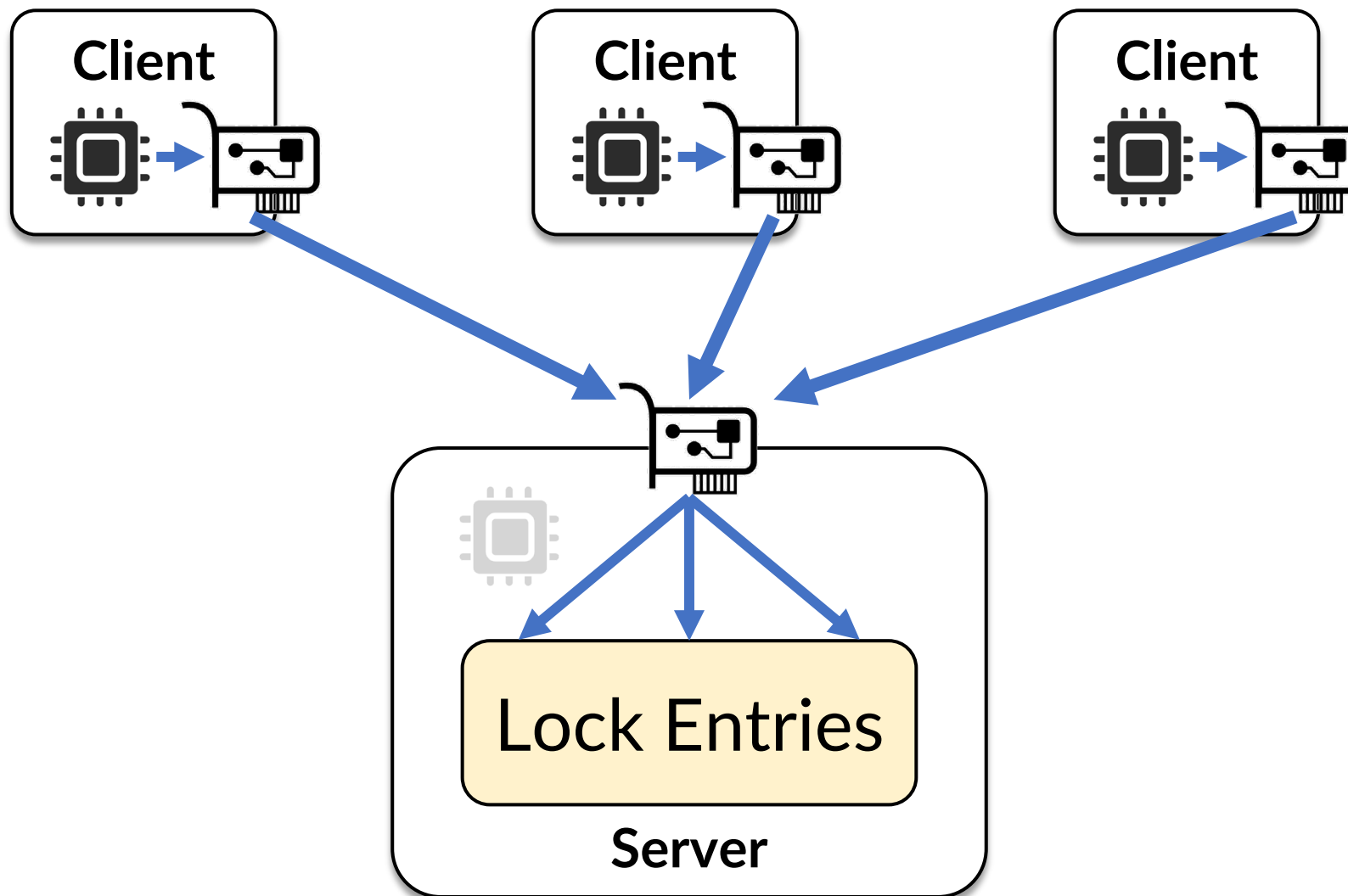
(Extended) Atomic verbs can ...

- Operate on larger-than-8B values
 - *16-byte atomics* with similar performance to original 8-byte atomics
- Compare-swap (CAS) on arbitrary bits
- Fetch-add (FAA) on separate fields

RDMA atomics are more powerful than CPU atomics.



(One-sided) RDMA Locks





RDMA Locks are Good

Because of ...

- Low latency
 - *Acquire/release a lock in an RDMA RTT (2~3 us)*
- High throughput
 - *370M packets/s for NVIDIA ConnectX-7 RNIC*
- Low CPU utilization
 - *Save power, no more CPU bottlenecks*



... But Problems Remain

One-sided = Server-inactive

- Client must **retry on failure**
 - *Retry traffic all heads to server*
- Retries are expensive
 - *RNIC queues up conflicting requests*

⇒ **Low performance under contention**

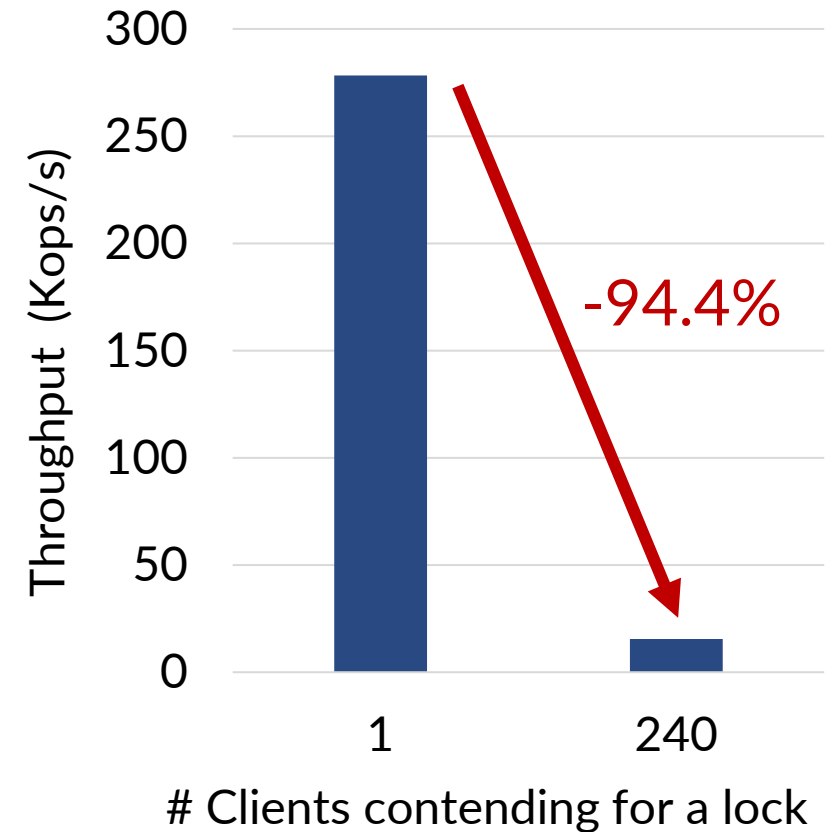


... But Problems Remain

One-sided = Server-inactive

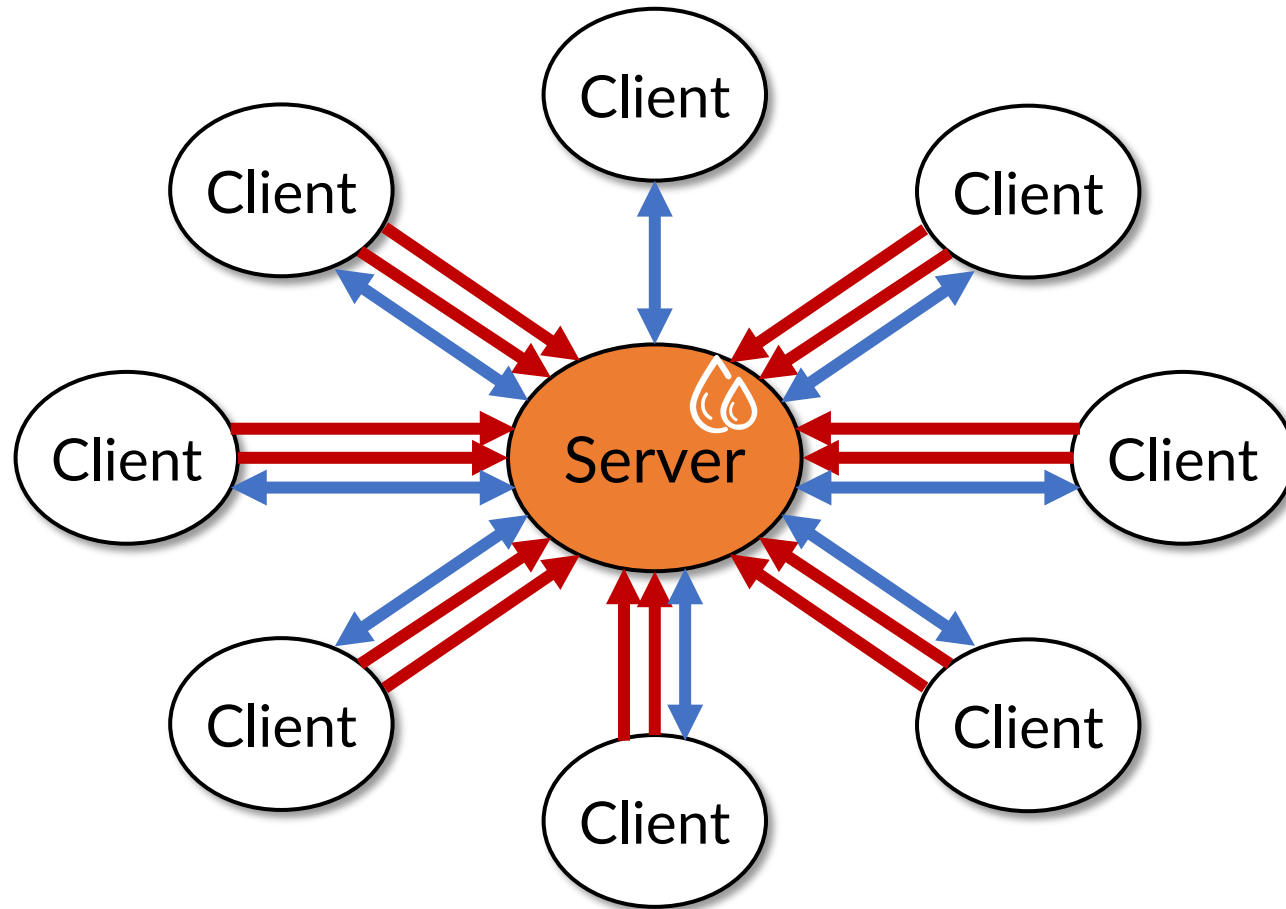
- Client must **retry on failure**
 - *Retry traffic all heads to server*
- Retries are expensive
 - *RNIC queues up conflicting requests*

⇒ **Low performance under contention**





Problem - Too Many Retries

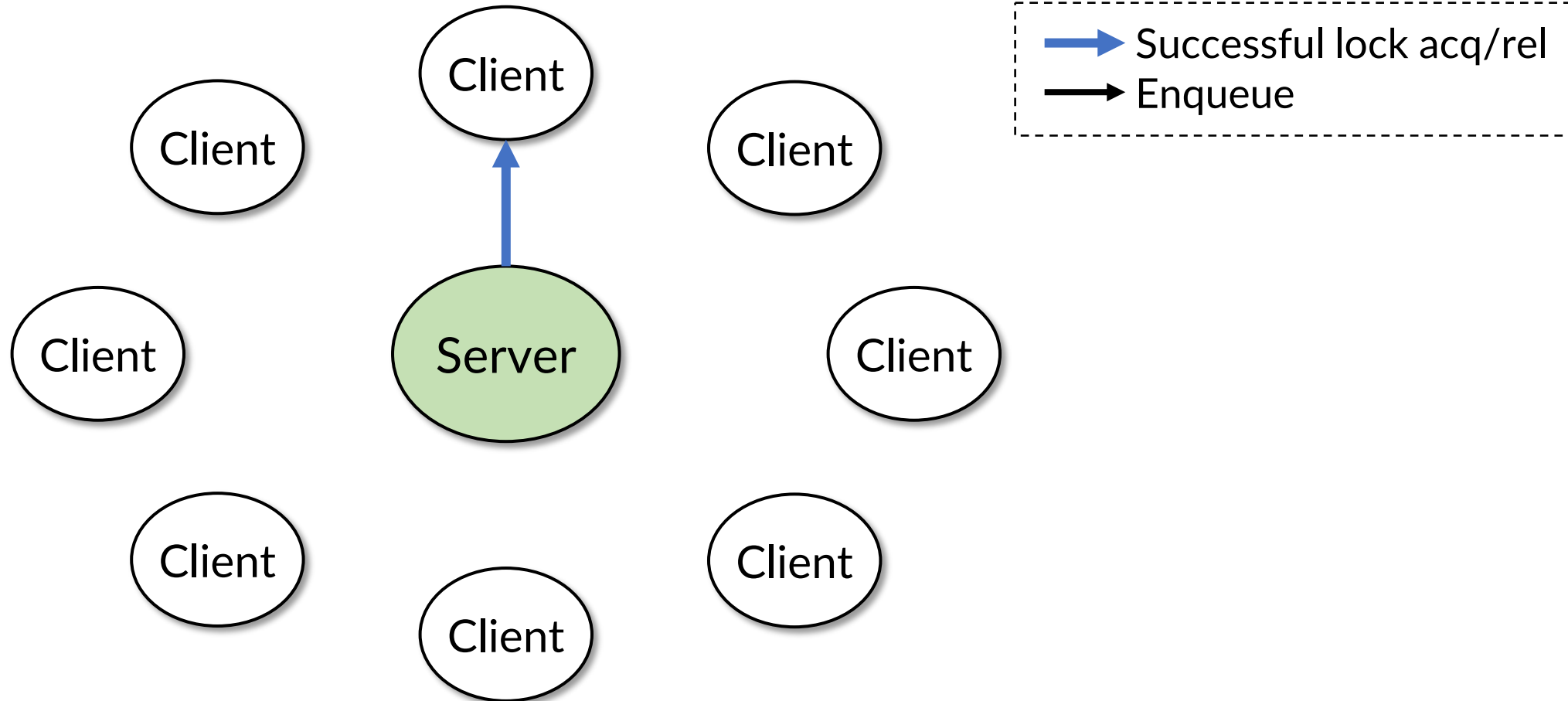


$\geq 85\%$
failed retries
in a typical CAS lock
(240 clients, *backoff*^[1] enabled)

[1] Ren et al. Scaling Up Memory Disaggregated Applications with SMART. ASPLOS (2024).

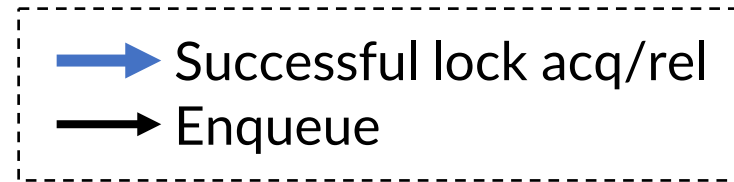
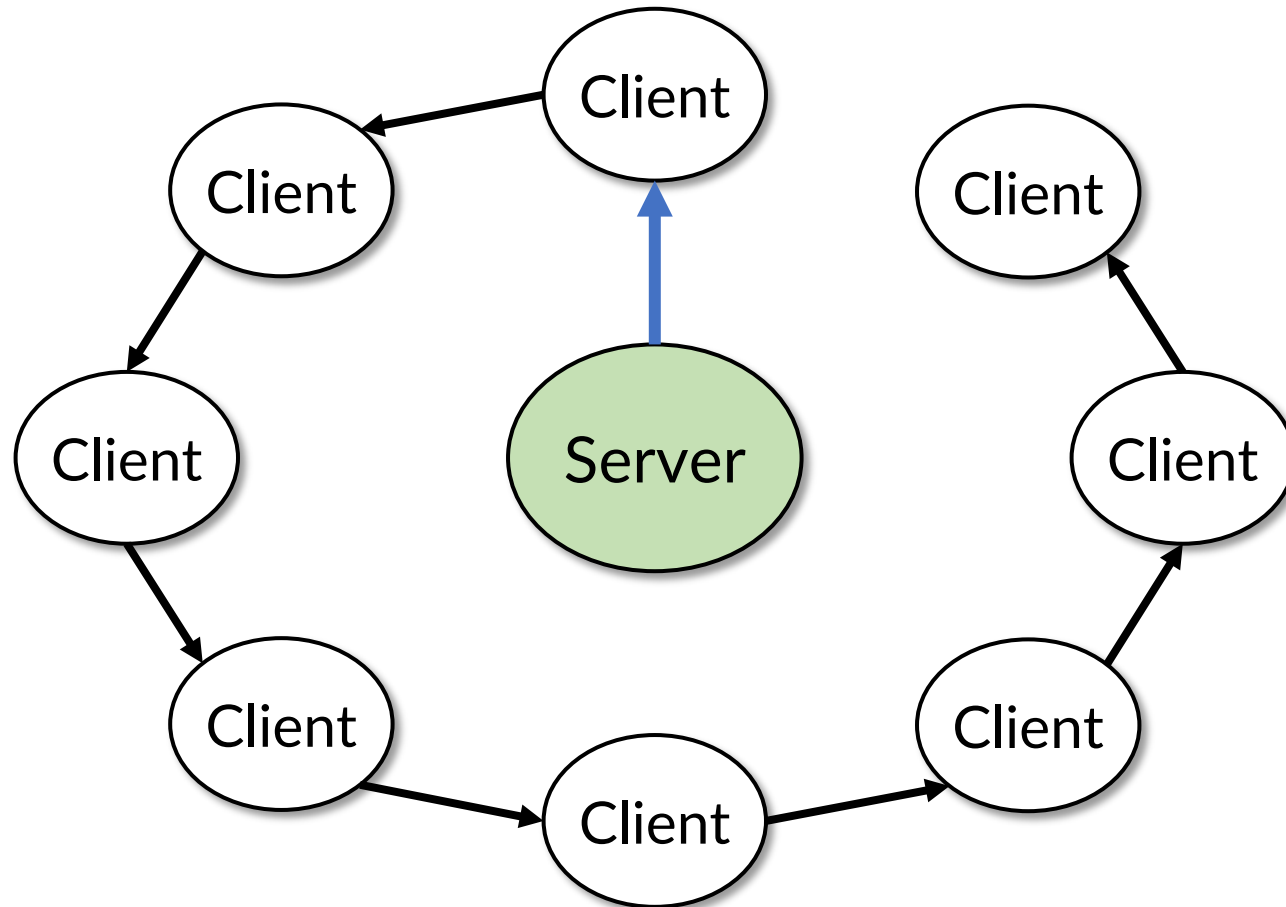


Solution - Let Clients **Hand over** Locks





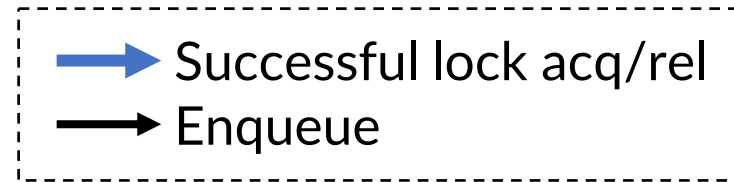
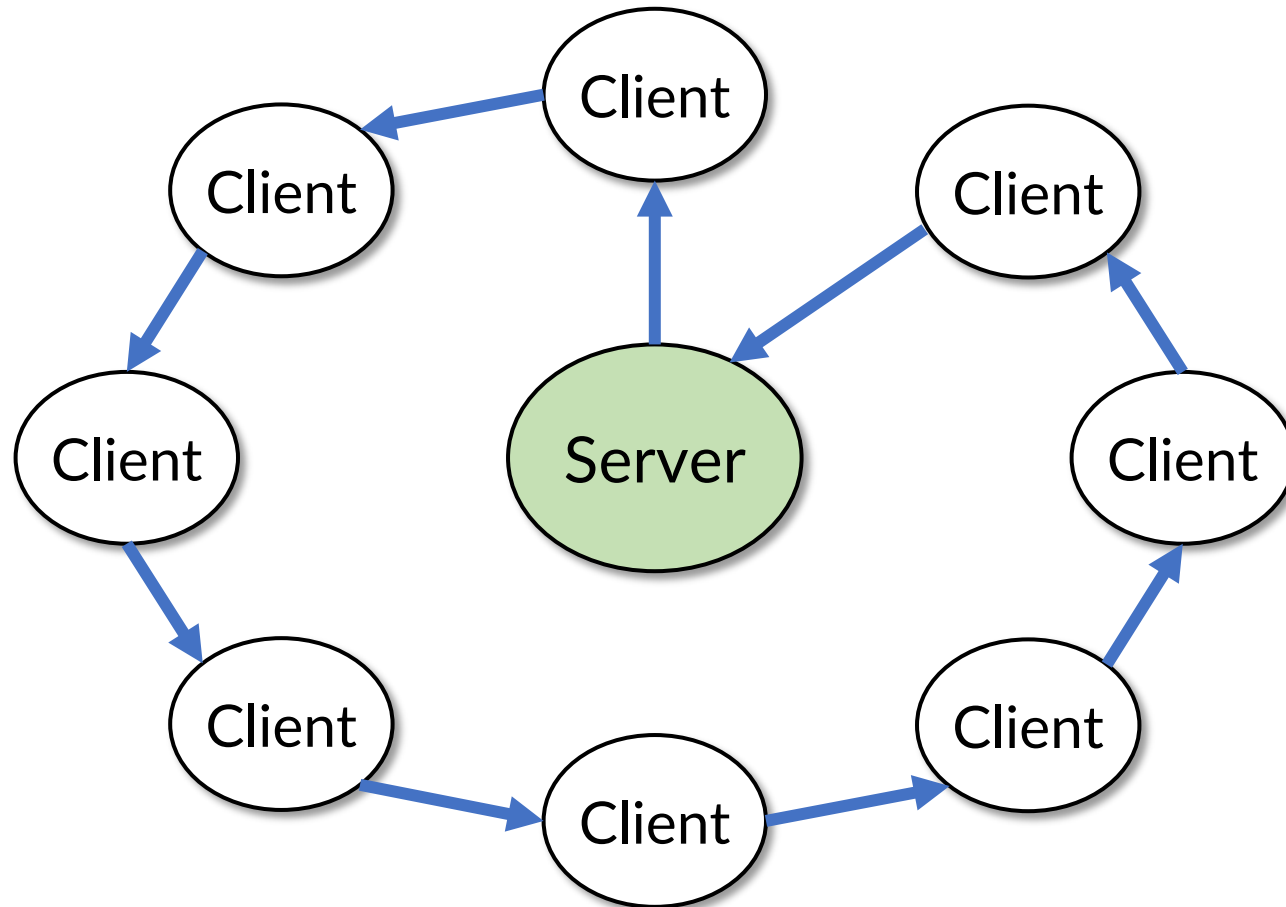
Solution - Let Clients **Hand over** Locks



- Clients **queue up**
 - Similarly to **MCS / CLH locks**
 - Cli-to-cli **direct messaging**



Solution - Let Clients **Hand over** Locks



- Clients **queue up**
 - Similarly to **MCS / CLH locks**
 - Cli-to-cli **direct messaging**
- Lock travels along the queue



Challenges

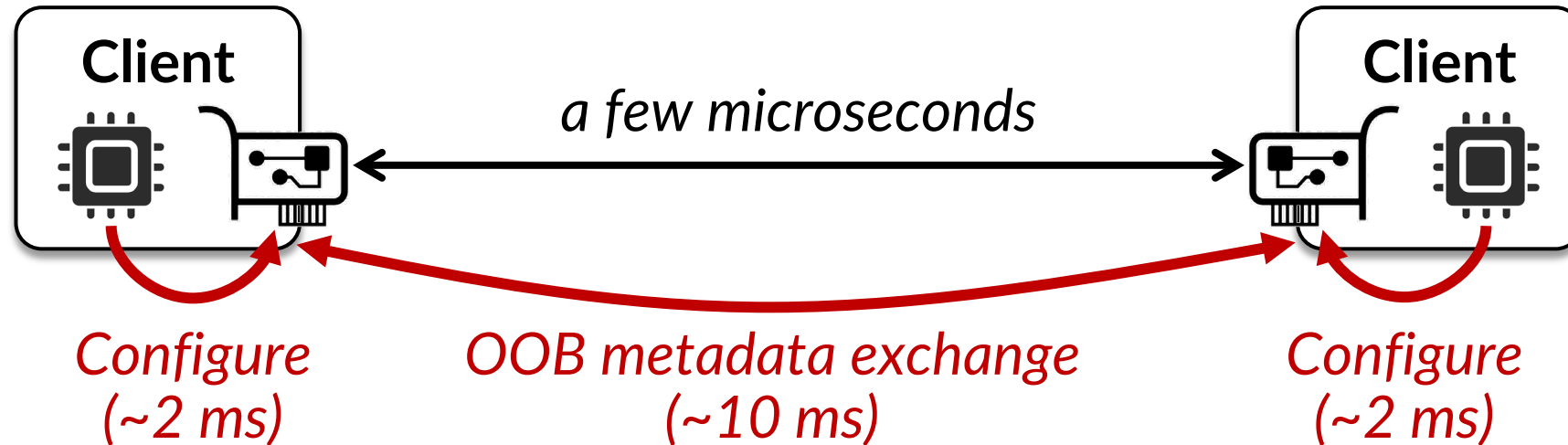
MCS lock^[1] is 34 yrs old. Why don't many RDMA locks use it?

- Expensive RDMA control path
 - *Difficult to enable cli-to-cli messaging*
- Long network RTT
 - *Single-node rwlock protocols won't work*
- Need of recovery
 - *Clients may fail*

[1] J.M. Mellor-Crummey. Algorithms for Scalable Synchronization on Shared-Memory Multiprocessors. ACM TOCS (1991).



Challenge #1: Expensive Control Path

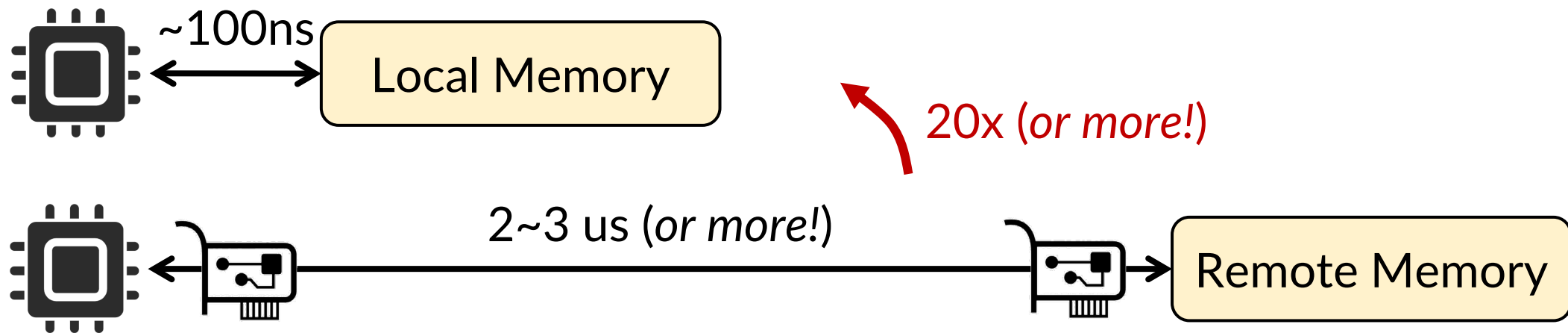


How to enable cli-to-cli messaging?

- All-to-all connection? *Unscalable*
- Connect at runtime? *Slow*



Challenge #2: Long RTT

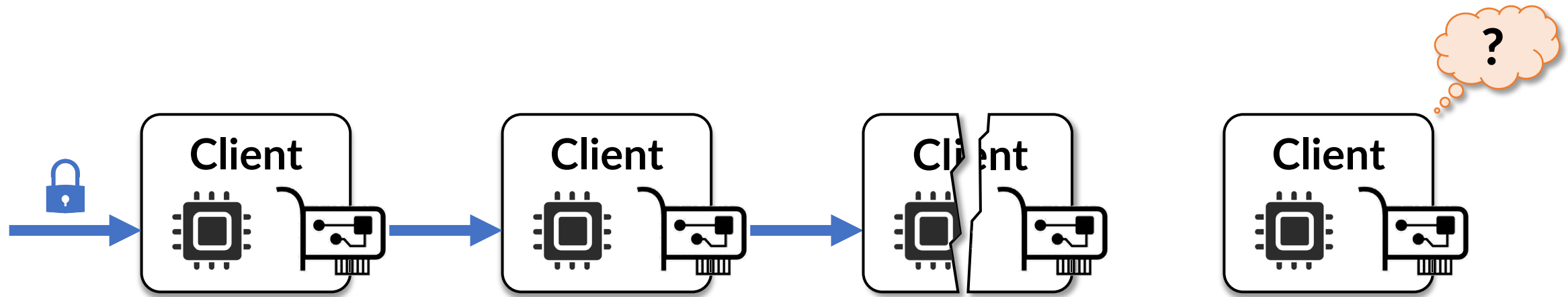


What lock semantics to implement?

- Mutex? *Low concurrency*
- Reader-writer lock? *Extra remote accesses*



Challenge #3: Fault Tolerance



How to tolerate client failures?

- How to detect a failure?
- How to correctly recover?



ShiftLock Outline

Idea of MCS lock, plus ...

- Scalable cli-to-cli communication
- Single-RTT rwlock protocol
- Fault tolerance

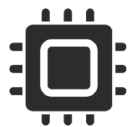


Basis - MCS Mutex

Lock Entry
(a pointer)

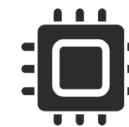
Tail null

ok false	next null
--------------------	---------------------



Client A

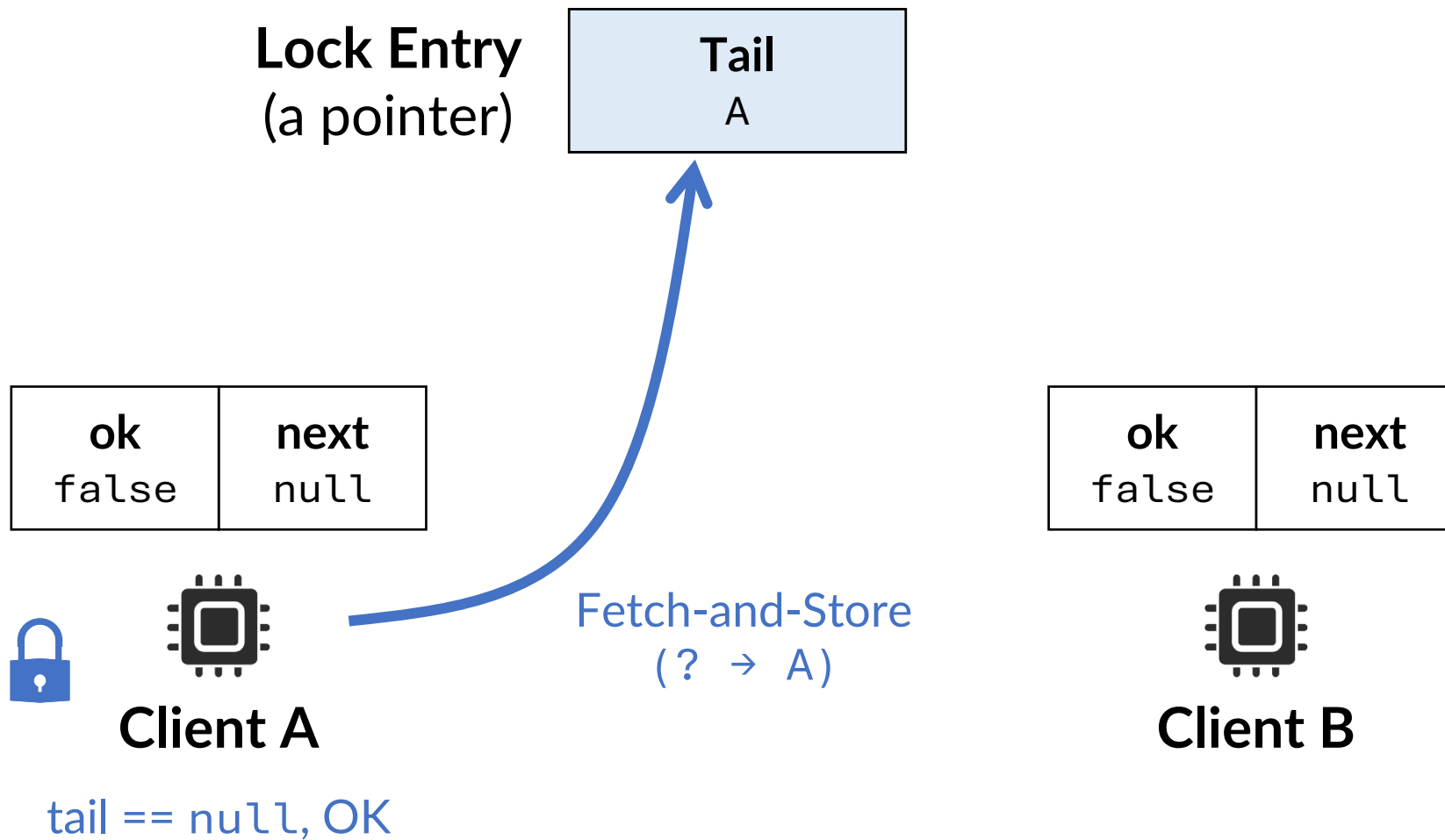
ok false	next null
--------------------	---------------------



Client B

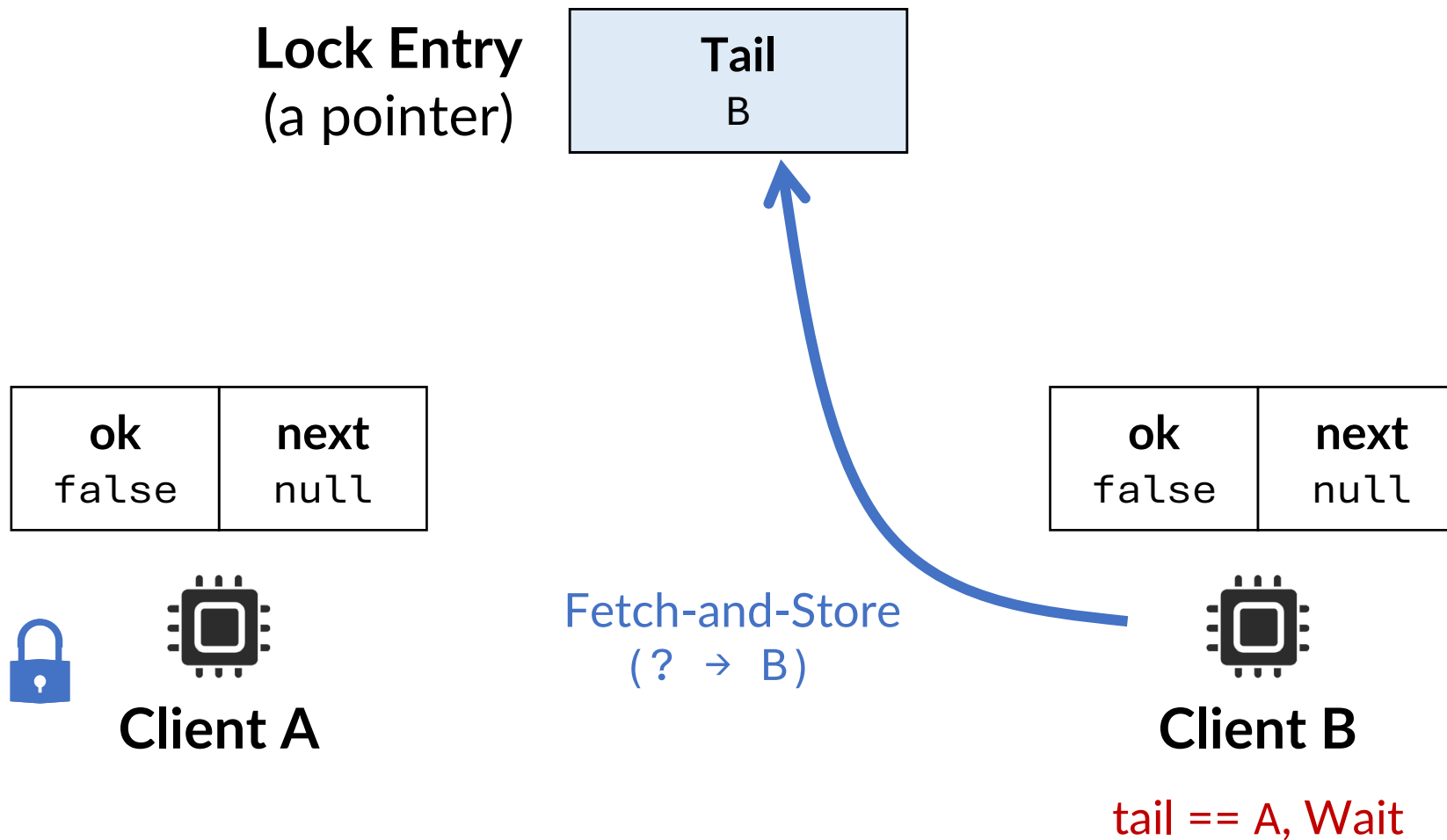


Basis - MCS Mutex





Basis - MCS Mutex

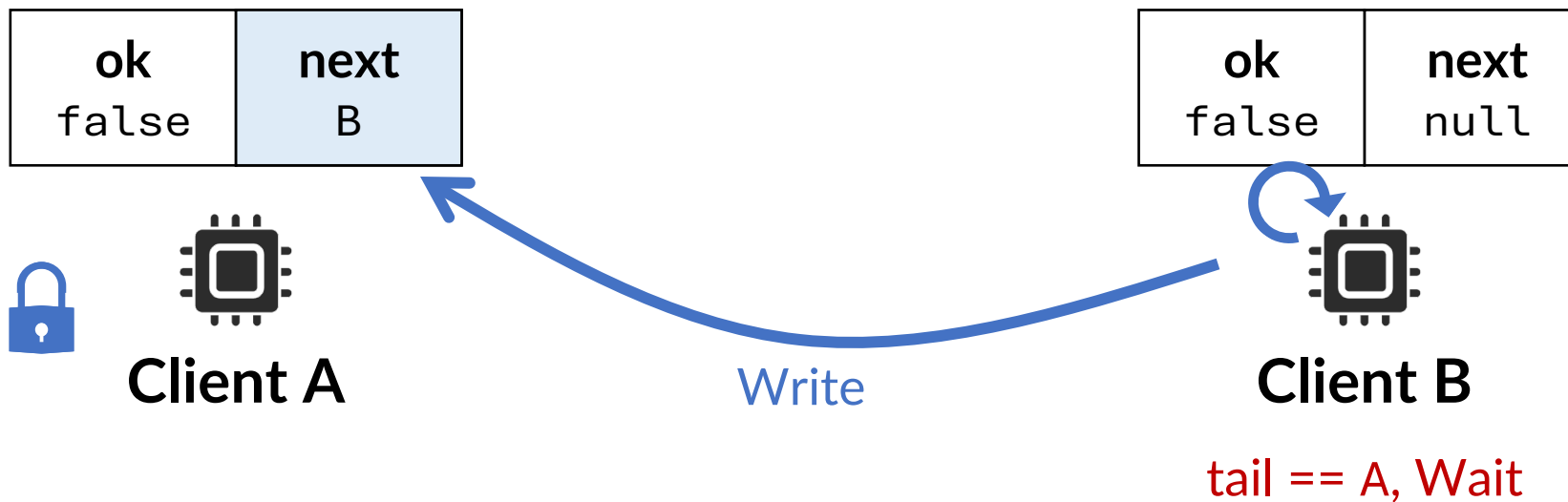




Basis - MCS Mutex

Lock Entry
(a pointer)

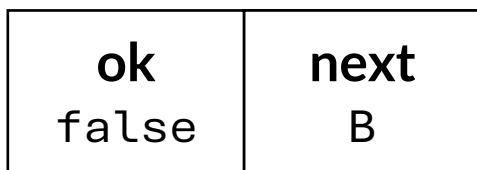
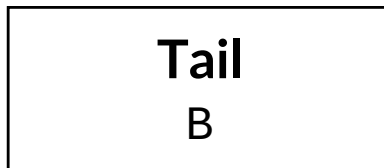
Tail B



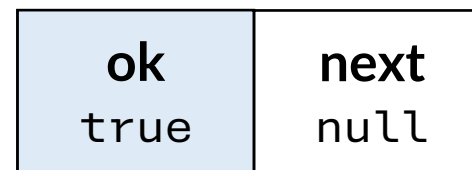


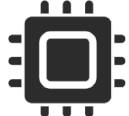

Basis - MCS Mutex

Lock Entry
(a pointer)




Client A
Unlock




Client B 



Write
Handover



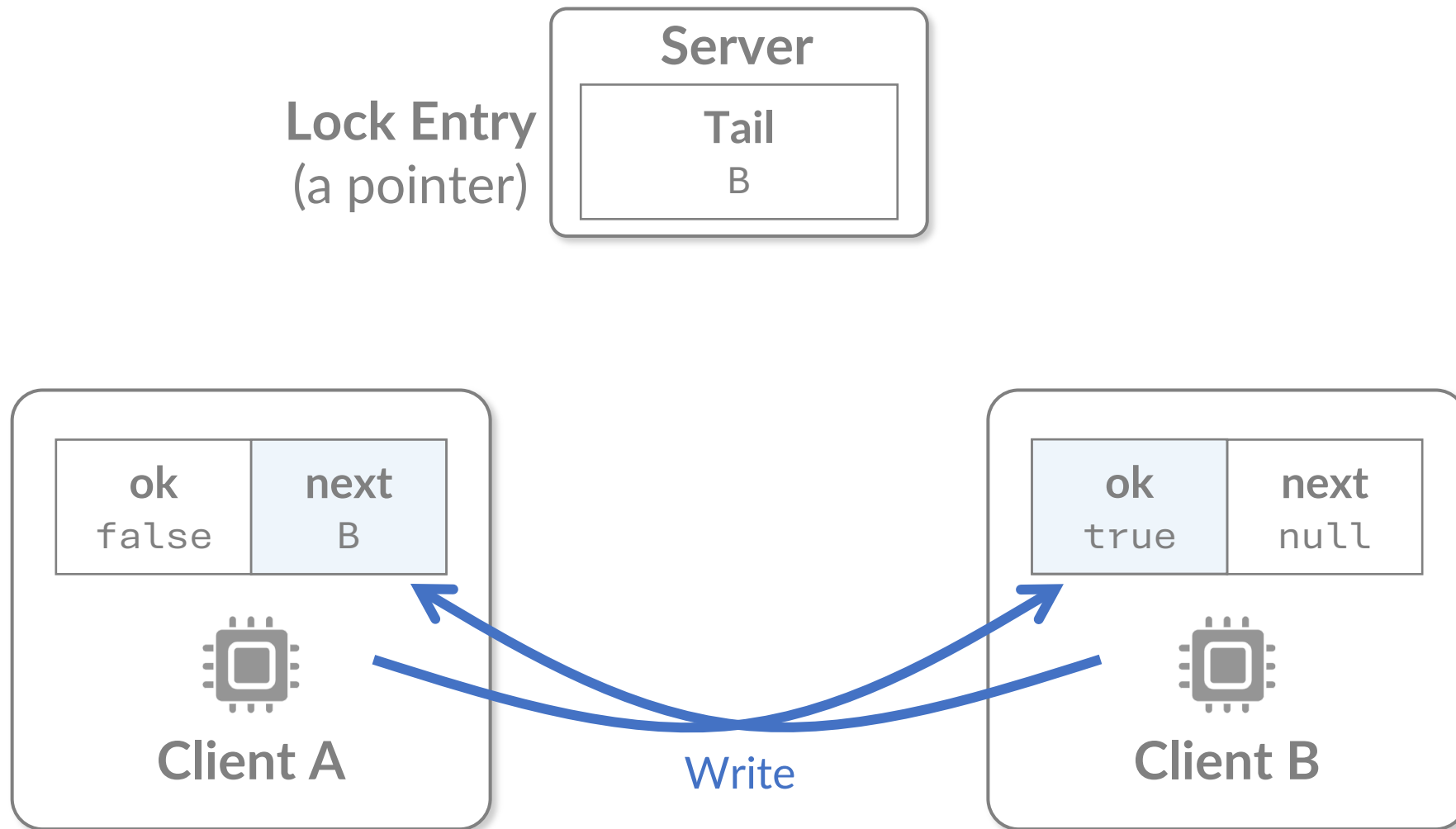
ShiftLock Outline

Idea of MCS lock, plus ...

- Scalable cli-to-cli communication
- Single-RTT rwlock protocol
- Fault tolerance

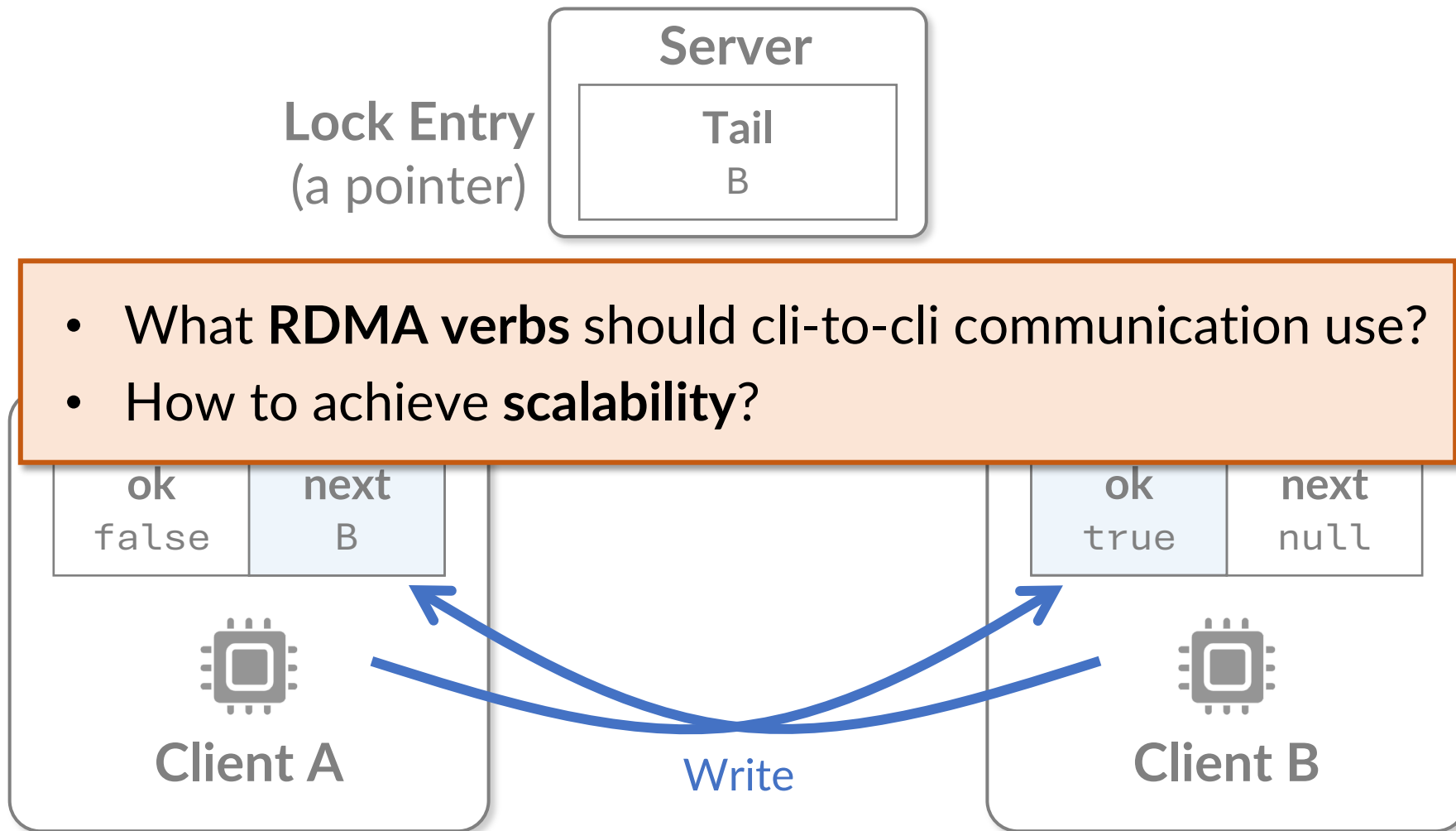


Design #1 - Cli-to-cli Communication





Design #1 - Cli-to-cli Communication





Design #1.1 - Using Two-sided Verbs

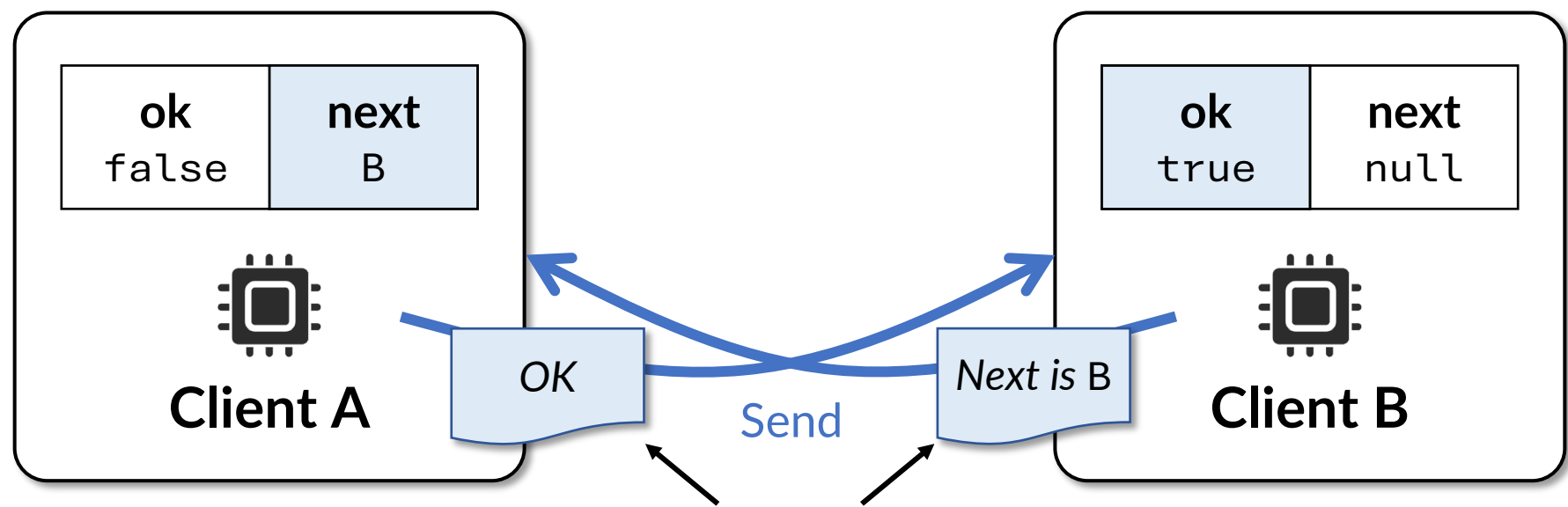
One-sided RDMA Write:

GID	LID	QPN	+	VAddr	RKey
128b	16b	24b		64b	32b

Two-sided RDMA Send:

GID	LID	QPN
128b	16b	24b

Save 96 bits of metadata

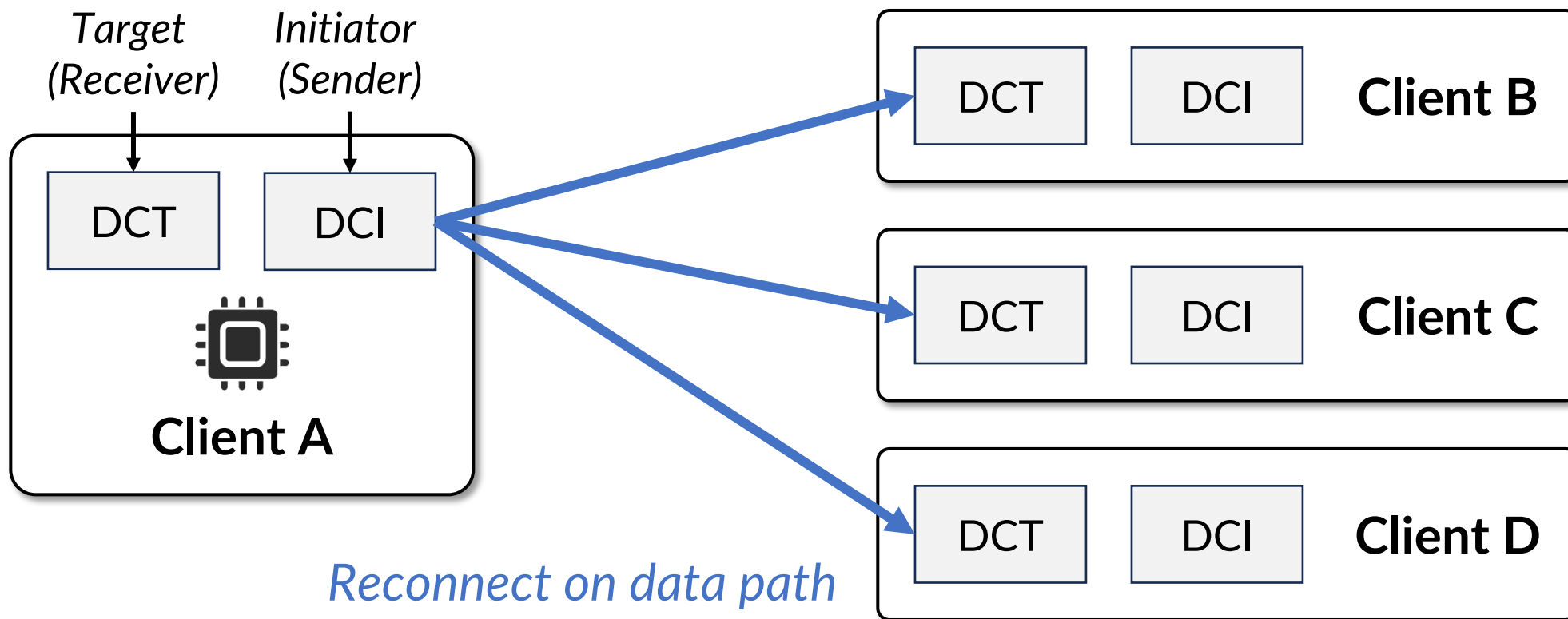


Processed only when needed



Design #1.2 - Scalability

RDMA Dynamic Connection



*Reconnect on data path
Overhead < 1 μs*



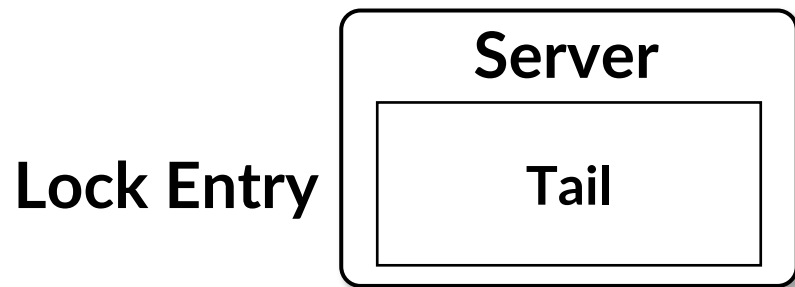
ShiftLock Outline

Idea of MCS lock, plus ...

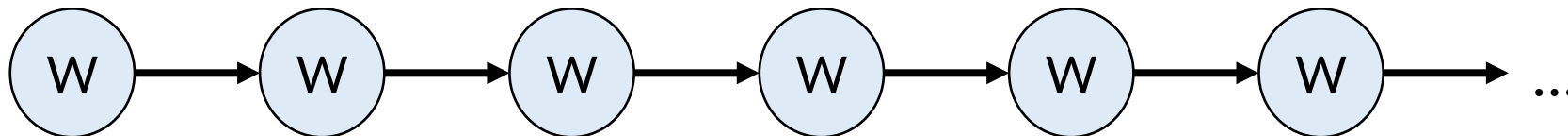
- Scalable cli-to-cli communication
- **Single-RTT rwlock protocol**
- Fault tolerance



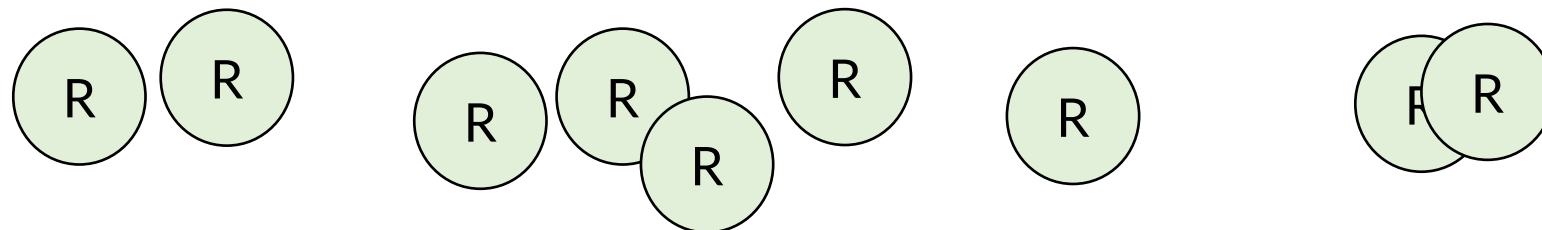
Design #2 - Rwlock Semantics



Writers
(in a queue)

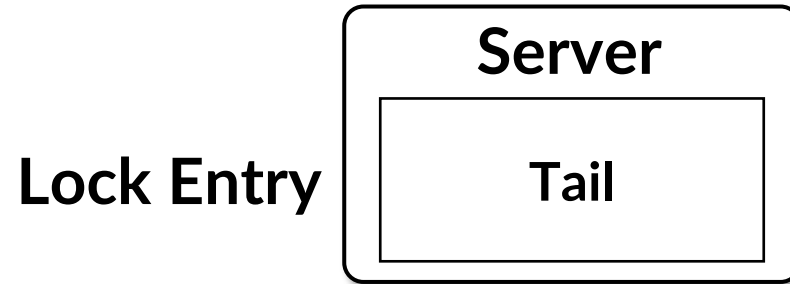


Readers?





Design #2 - Rwlock Semantics

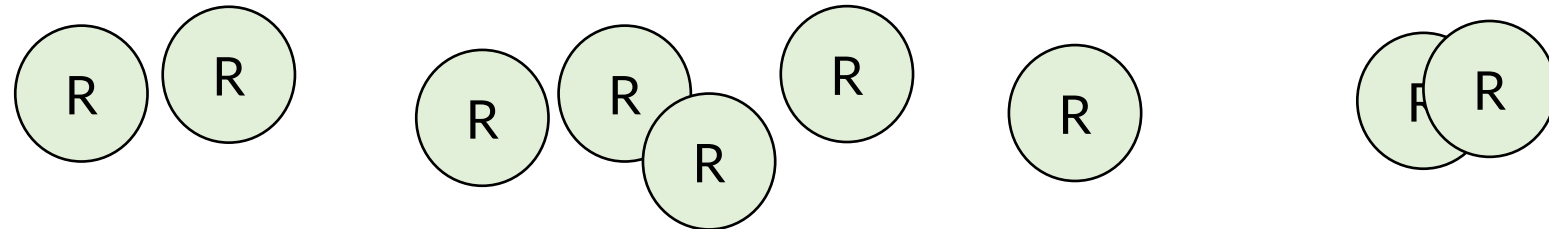


- How to maintain readers?
- How to transfer lock between readers & writers?

Writers
(in a queue)

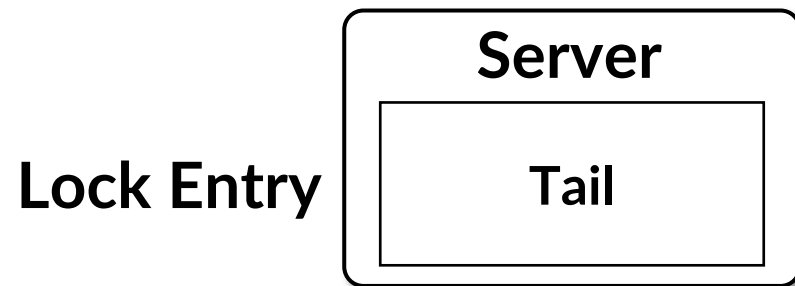


Readers?

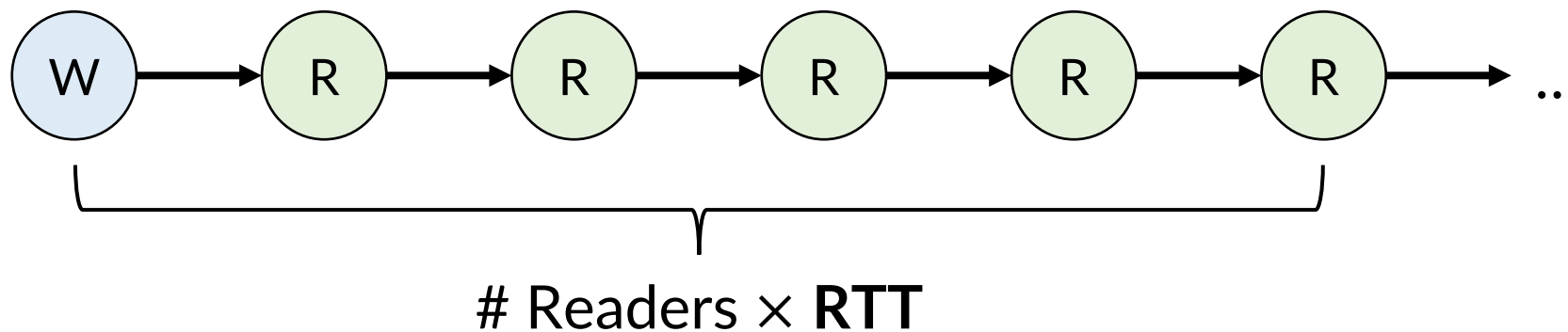




Design #2.0 - Queues Unfit for Readers



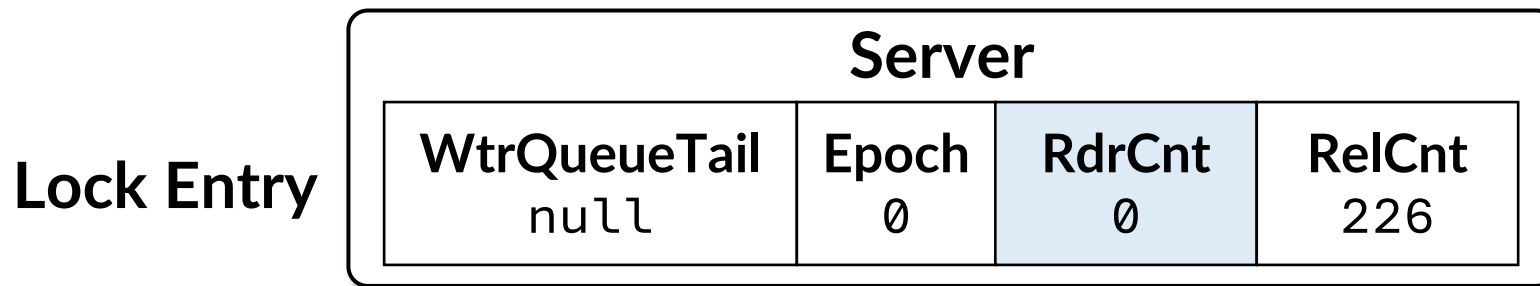
A unified queue?



Latency too high, NG!



Design #2.1 - Reader Counter



Count current readers

- +1 when try to acquire
- -1 when release

Writer acquire:
(except handover)

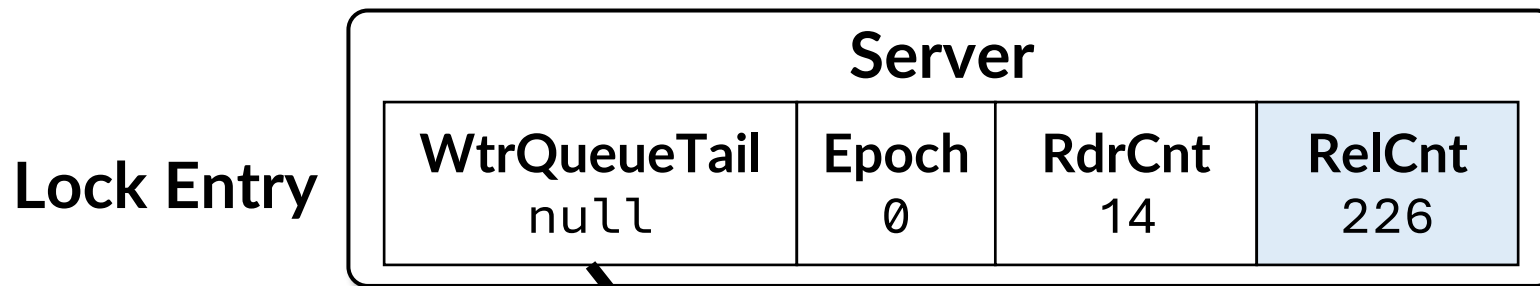
$\{ \text{WtrQueueTail} == \text{null} \}$ && $\{ \text{RdrCnt} == 0 \}$
(no writers) *(no readers)*

Reader acquire:

$\{ \text{WtrQueueTail} == \text{null} \}$
(no writers)



Design #2.2 - RLock \rightarrow WLock



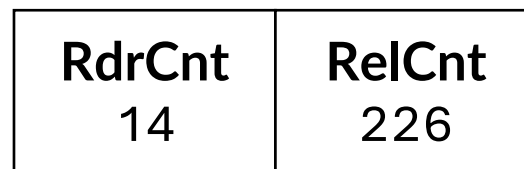
Count lock releases

- +1 upon every release

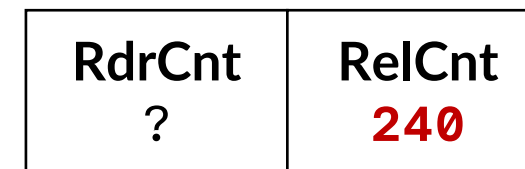
When enqueueing,
atomically read
the entire entry

Readers before have all unlocked,
acquire!

Writer acquire:
(when reader holds the lock)

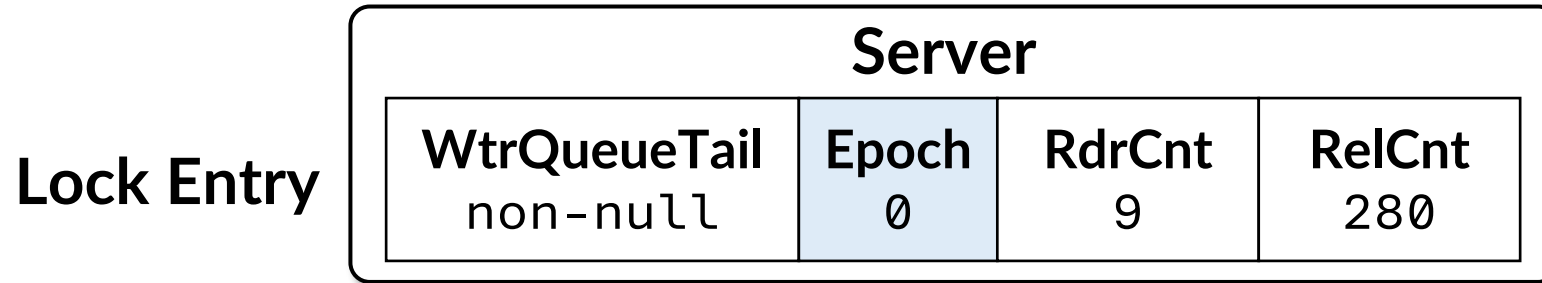


Wait
→



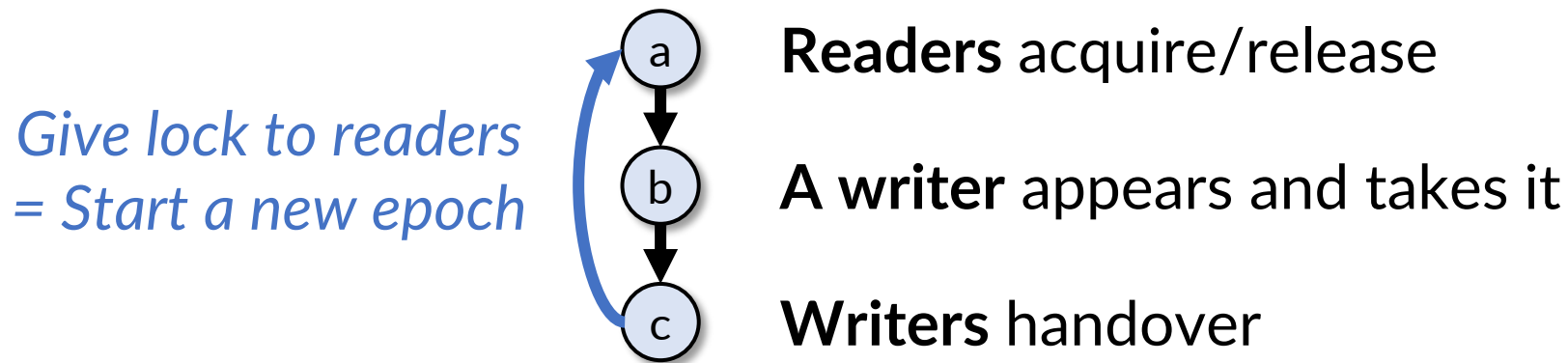


Design #2.3 - WLock \rightarrow RLock



Give lock to readers

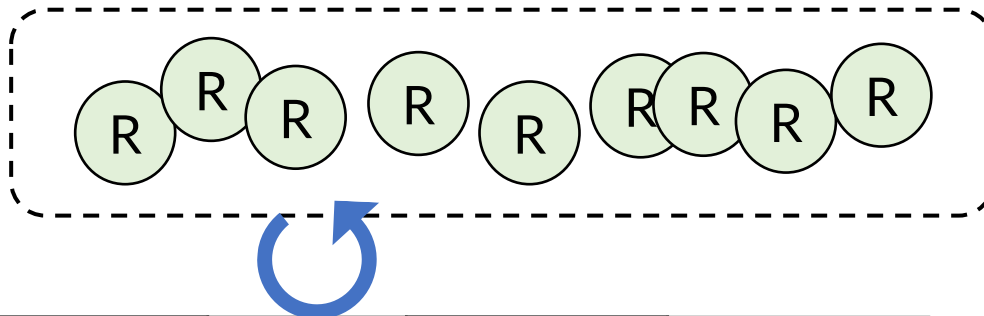
Lock lifetime contains many epochs. An epoch contains 3 steps:





Design #2.3 - WLock \rightarrow RLock

Reader acquire:
(when writer holds the lock)



WtrQueueTail non-null	Epoch 0	RdrCnt 9	RelCnt 280
---------------------------------	-------------------	--------------------	----------------------



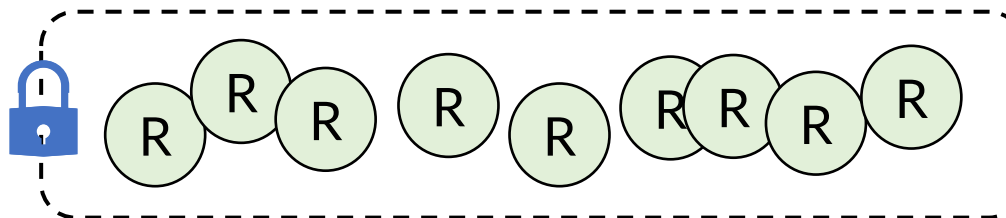
Writer A
next = B

Writer B



Design #2.3 - WLock \rightarrow RLock

Reader acquire:
(when writer holds the lock)



Epoch changed,
acquire

WtrQueueTail non-null	Epoch 1	RdrCnt 9	RelCnt 281
---------------------------------	-------------------	--------------------	----------------------

Writer A
next = B

Writer B

Flip

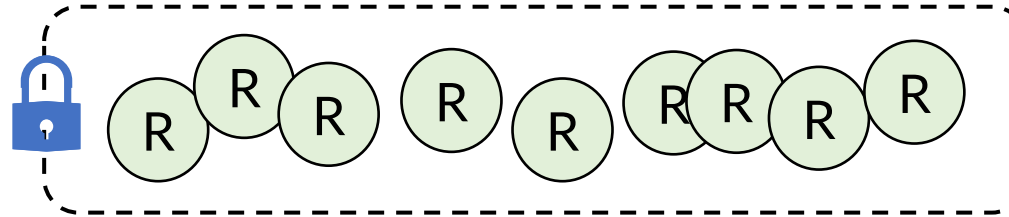
Increment

RdrCnt = 9, RelCnt = 280



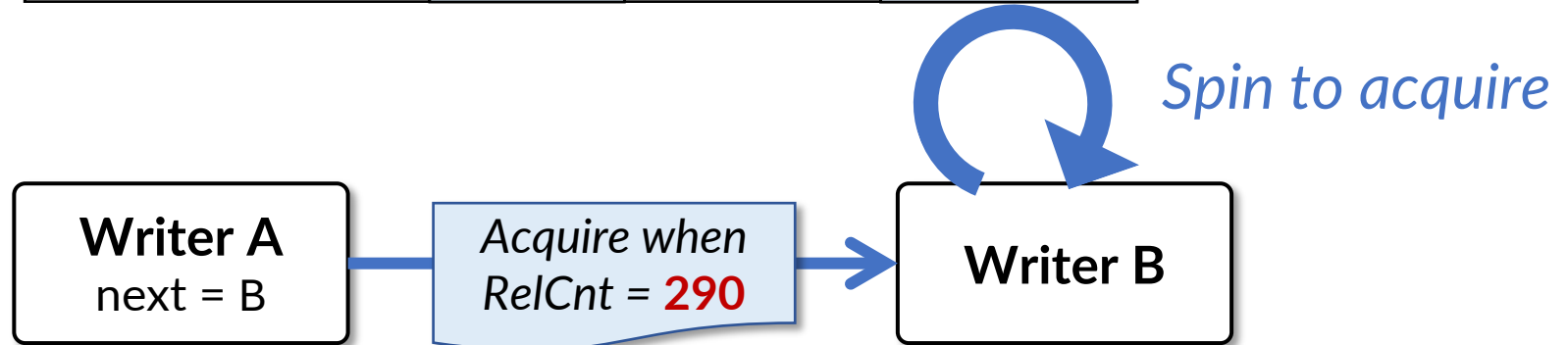
Design #2.3 - WLock \rightarrow RLock

Reader acquire:
(when writer holds the lock)



Epoch changed,
acquire

WtrQueueTail non-null	Epoch 1	RdrCnt 9	RelCnt 281
---------------------------------	-------------------	--------------------	----------------------



$$290 = (280 + 1) + 9$$



Design #2.3 - WLock \rightarrow RLock

Reader a
(when writer h

Readers acquire when Epoch **changes**
Epoch is not an “mode” that “must be correct”
Guarantees **single-RTT acquisition** w/o contention

ch changed,
acquire

wrtQueueTail Epoch RelCnt RelCnt

Writers **proactively** transfer the lock
Criteria: consecutive handovers \geq threshold

acquire

next = B

RelCnt = 290

Writer B

$$290 = (280 + 1) + 9$$



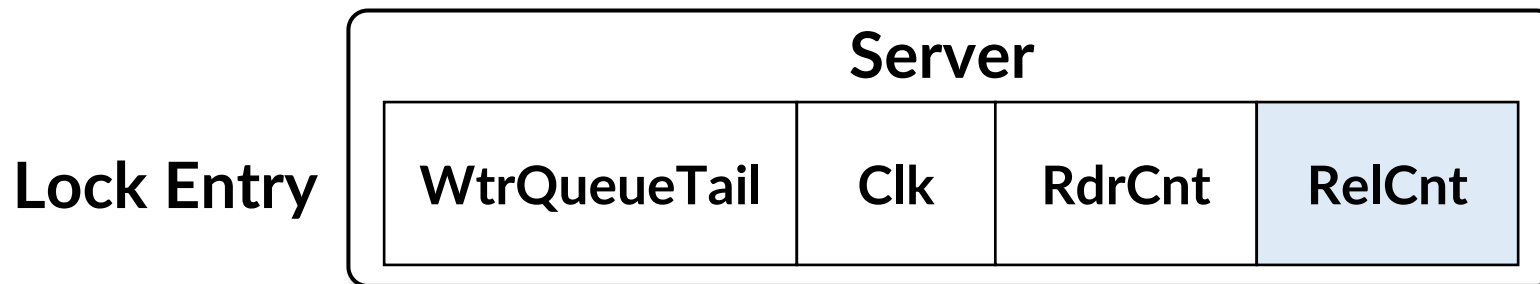
ShiftLock Outline

Idea of MCS lock, plus ...

- Scalable cli-to-cli communication
- Single-RTT rwlock protocol
- **Fault tolerance**

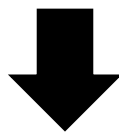


Design #3.1 - Detect Failures



Observation:

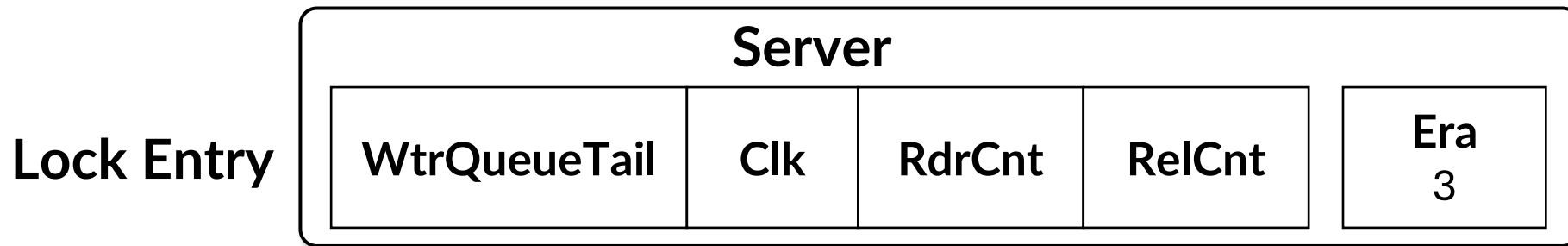
Client failure = eventually, no one will *release* the lock
= *RelCnt will stop changing*



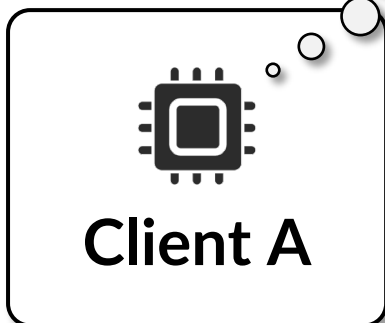
Detect a failure if *RelCnt keeps unchanged* for a certain period
(i.e., leasing)



Design #3.2 - Two-sided Recovery



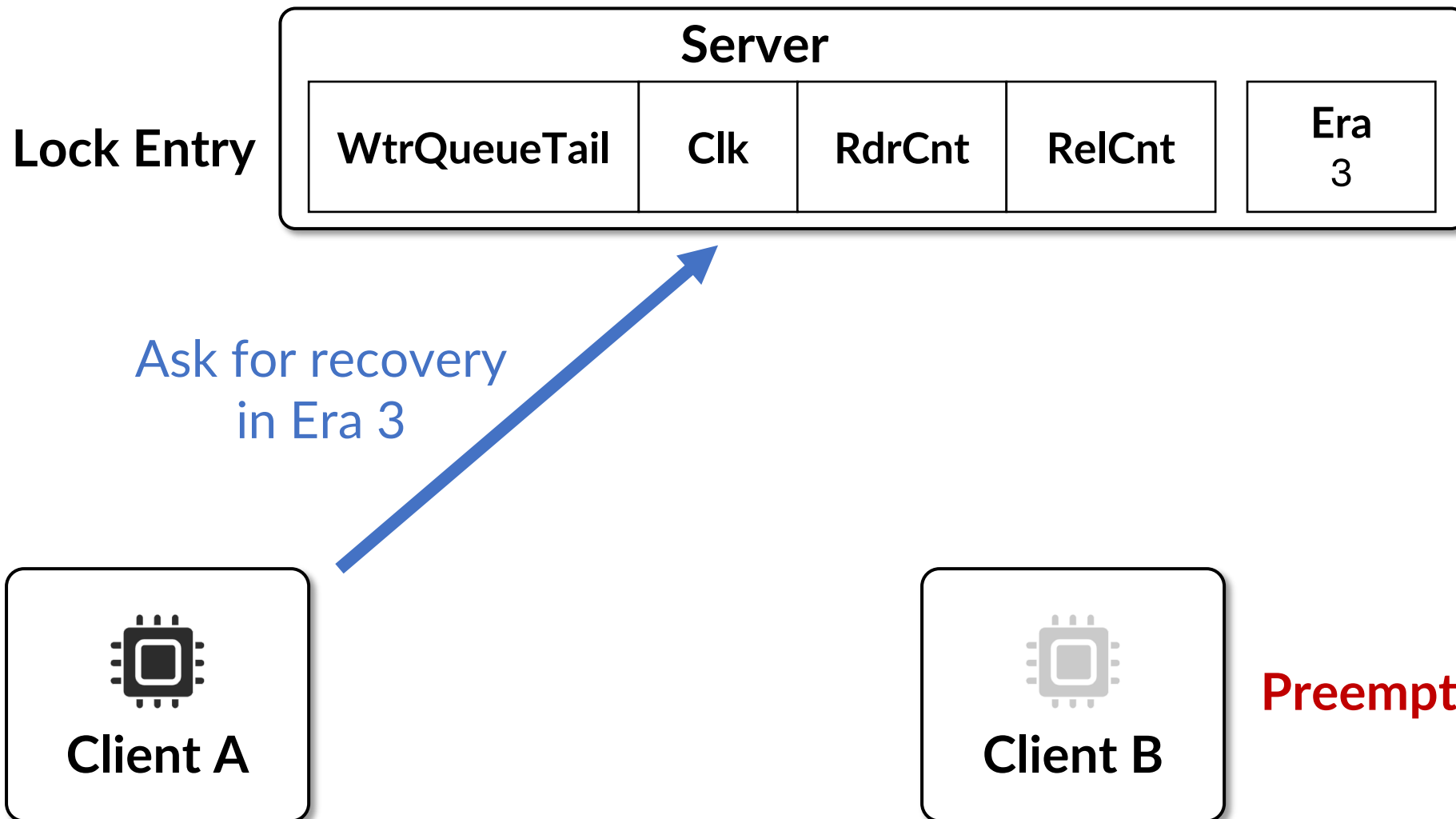
*Failure detected,
I will recover!*



Preempted by OS

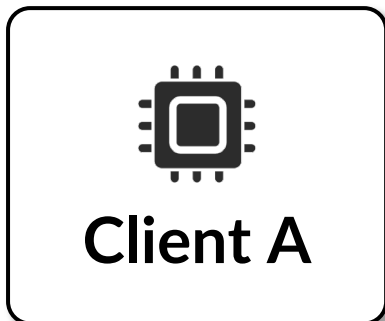
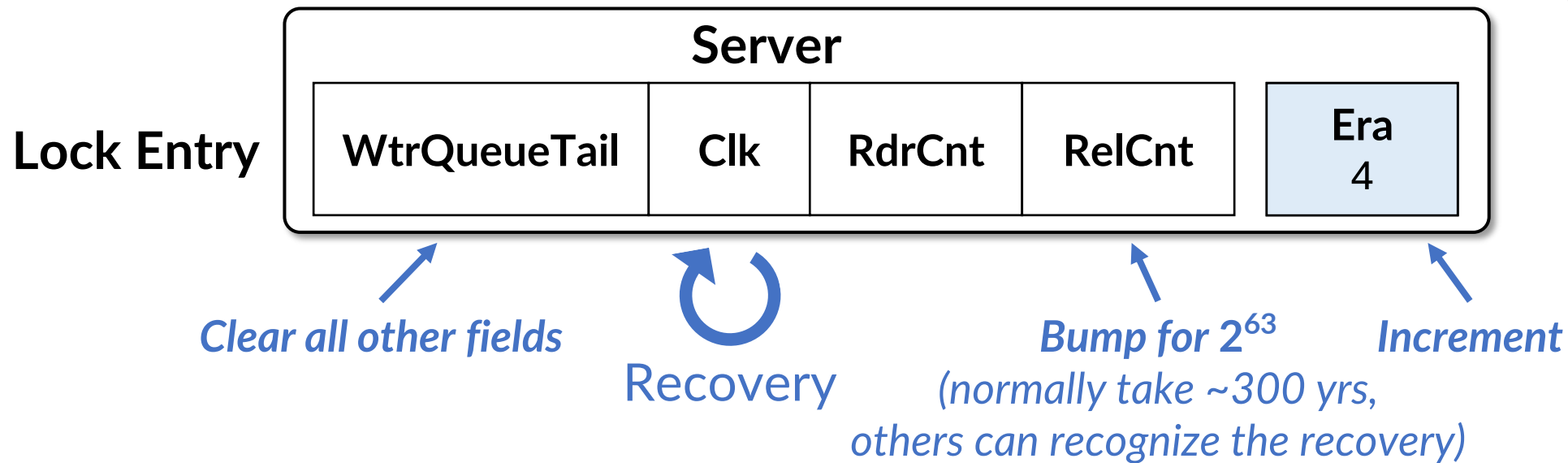


Design #3.2 - Two-sided Recovery





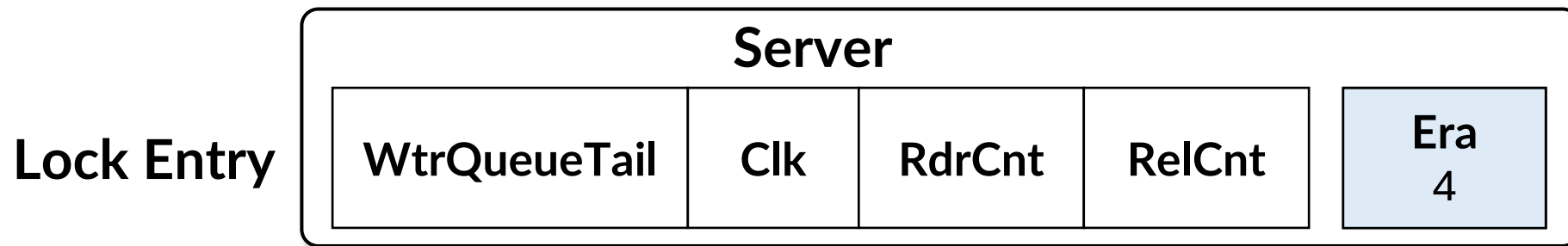
Design #3.2 - Two-sided Recovery



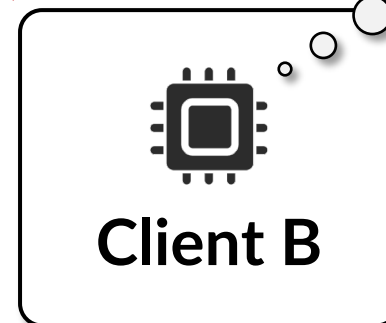
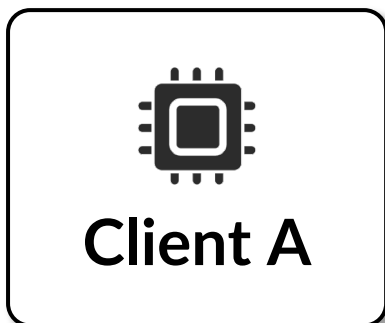
Preempted by OS



Design #3.2 - Two-sided Recovery



Ask for recovery in era 3
(will definitely fail)





Design Summary

- Scalable cli-to-cli communication
 - *Two-sided RDMA on Dynamic Connection*
- Single-RTT rwlock protocol
 - *Minimum latency w/o contention*
- Fault tolerance
 - *Correct two-sided recovery*

See our paper for more details! 😊



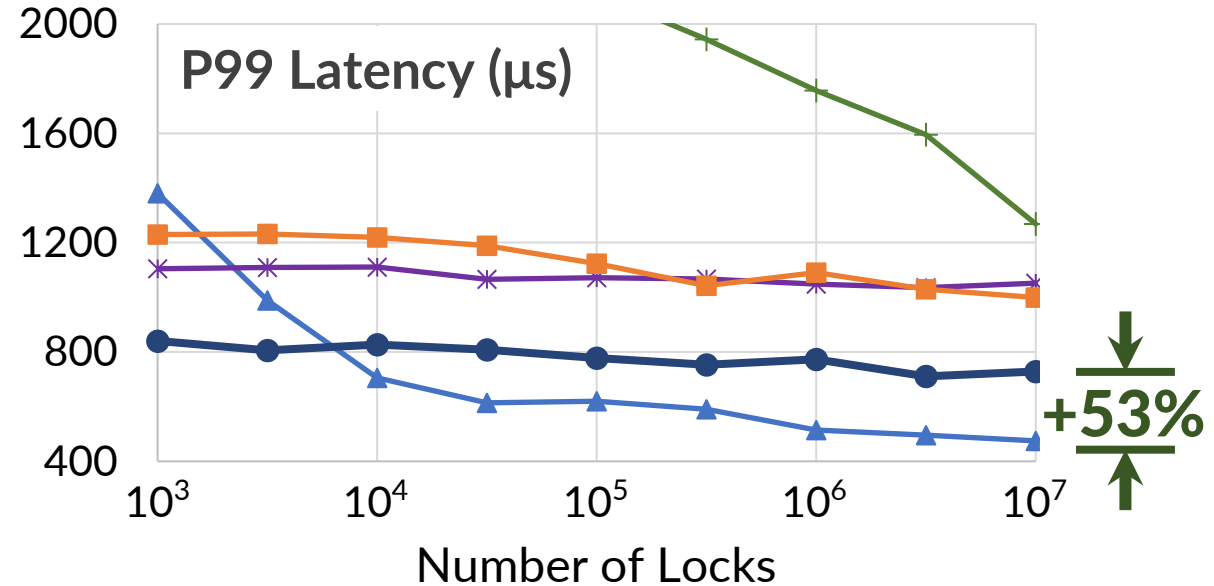
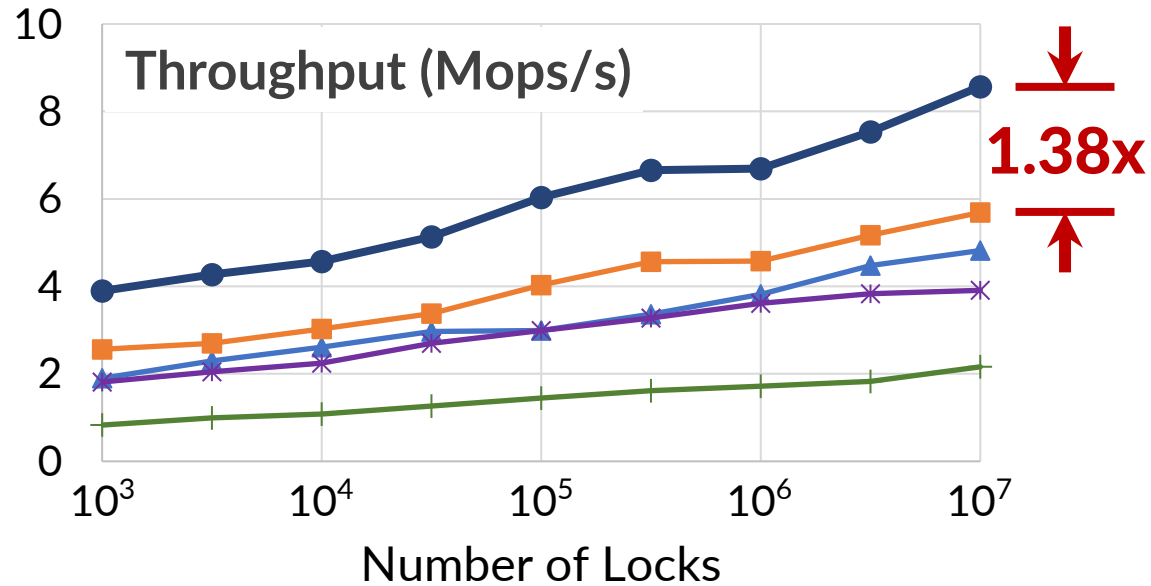
Evaluation

- Implemented in **Rust**, **~7.2k LoC**
- Compared with:
 - DSLR [SIGMOD'18], RMA-RW [HPDC'16], CAS (+Backoff), RPC, ...
- 1 server + 5 clients
- 2x Intel® Xeon® E5-2650v4 CPUs
 - $2 \times 24 = 48$ cores, $48 \times 5 = 240$ client threads
- Mellanox ConnectX-5 100Gbps RNIC + QM8790 switch



Overall Performance

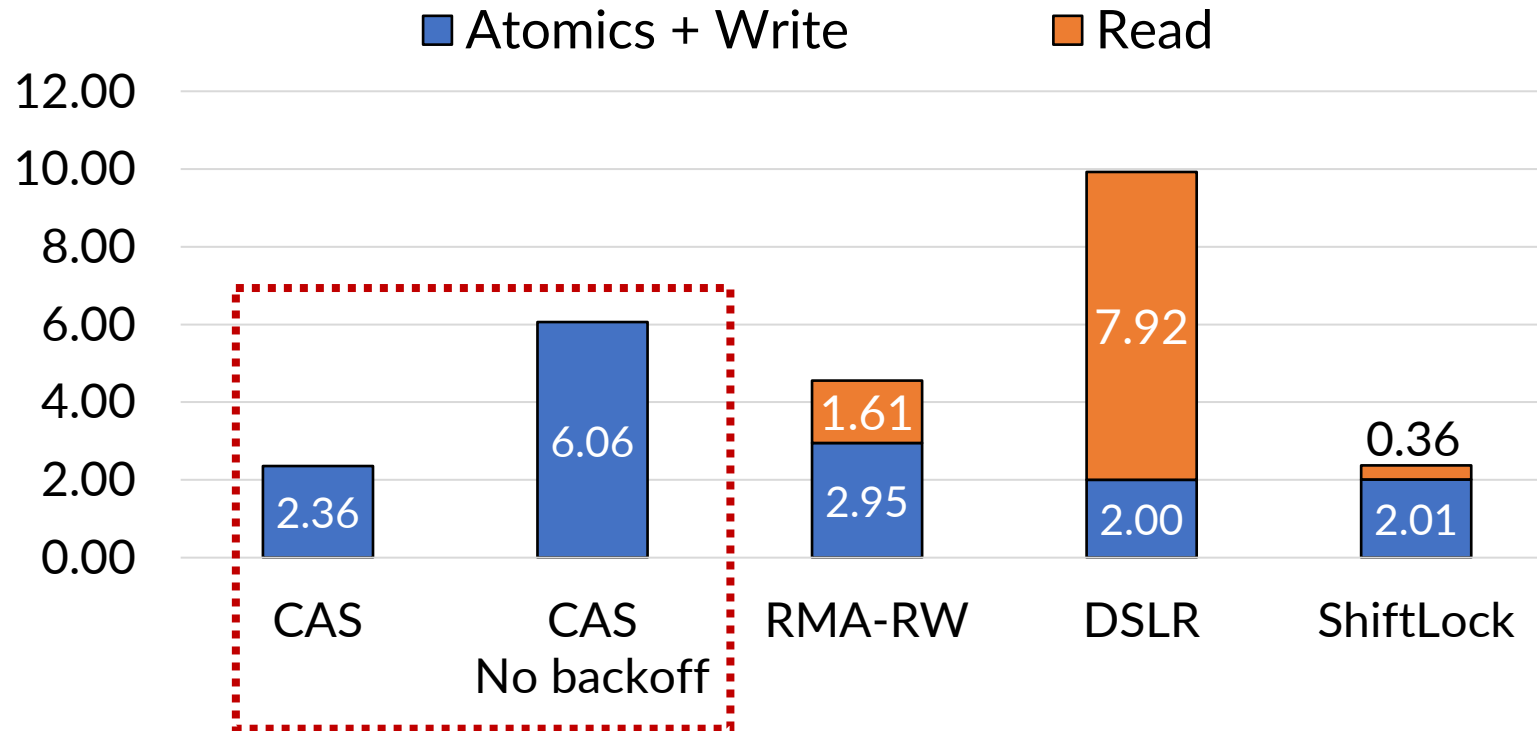
▲ CAS * DSLR + RPC ■ RMA-RW ● ShiftLock



- ShiftLock improves throughput & latency for rwlock at different scales
- Compared to mutex, writers wait longer (*due to actively giving locks to readers*)



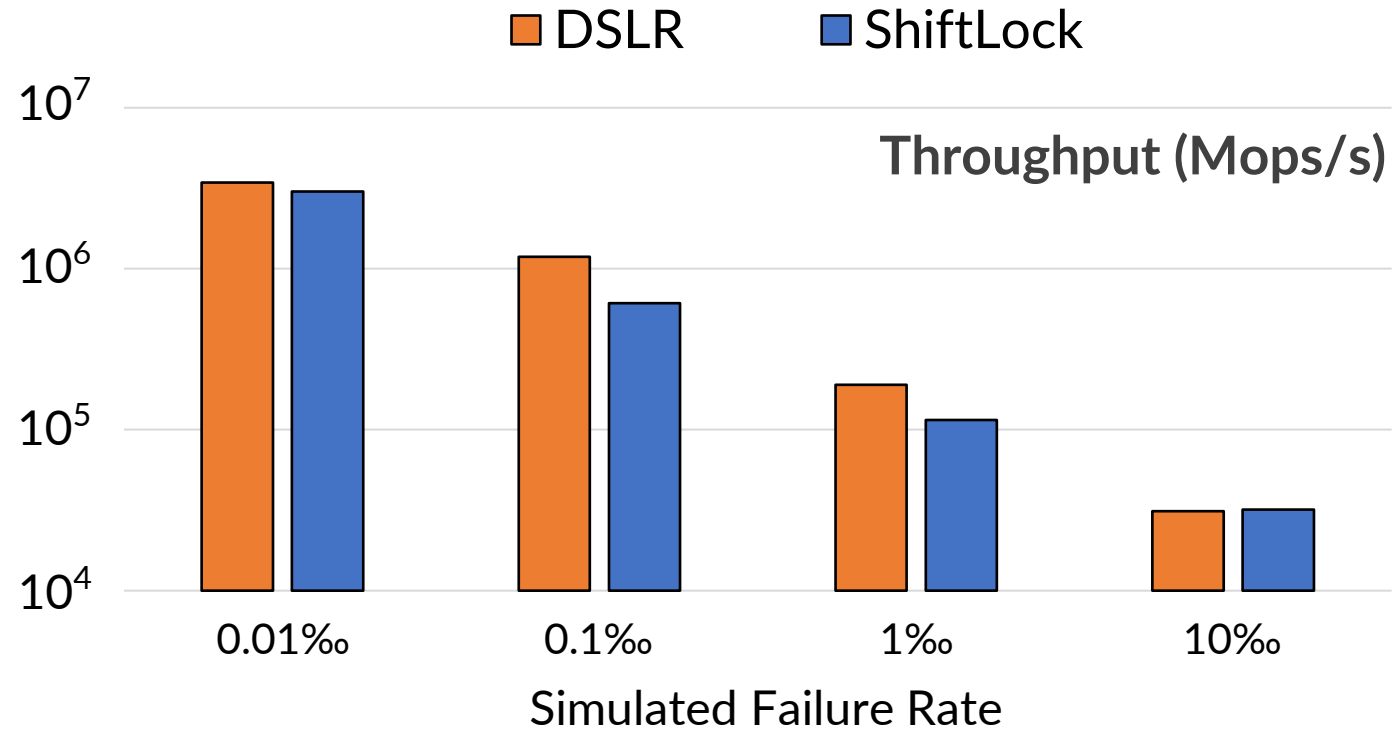
One-sided Verbs Received by the Server



- ShiftLock **effectively reduces retries** (*server RNIC has the lowest load*)
- Backoff has impressive effects, but **cannot solve the retry problem entirely**



Appendix: Recovery Performance



ShiftLock has **lower (but similar)** recovery performance than existing RDMA locks



Thank you! Questions are welcome.

- **ShiftLock**: Mitigate One-sided RDMA Lock Contention via Handover
- **MCS rwlock**, with ...
 - Scalability (*two-sided cli-to-cli messaging with Dynamic Connection*)
 - Low latency (*guaranteed single-RTT w/o contention*)
 - Fault tolerance (*correct recovery*)
- **Fewer retries**, higher throughput
- Artifact at <https://github.com/thustorage/shiftlock>
- Contact: gaoj20@mails.tsinghua.edu.cn

