

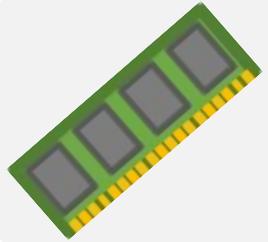
NrOS: Effective Replication and Sharing in an Operating System

Ankit Bhardwaj, Chinmay Kulkarni, Reto Achermann, Irina Calciu,
Sanidhya Kashyap, Ryan Stutsman, Amy Tai, and Gerd Zellweger



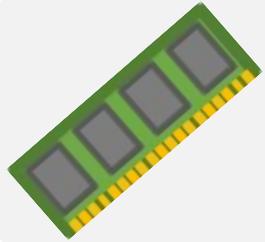
Kernel development is hard

Kernel development is hard



Manual memory management

Kernel development is hard

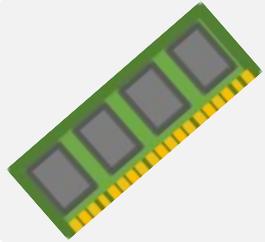


Manual memory management



Type-unsafe code

Kernel development is hard



Manual memory management



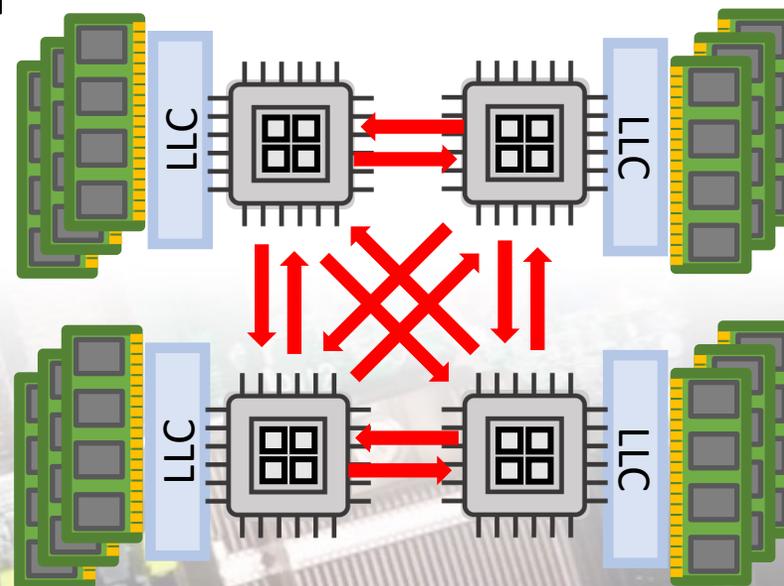
Type-unsafe code



Complex and asynchronous event handling

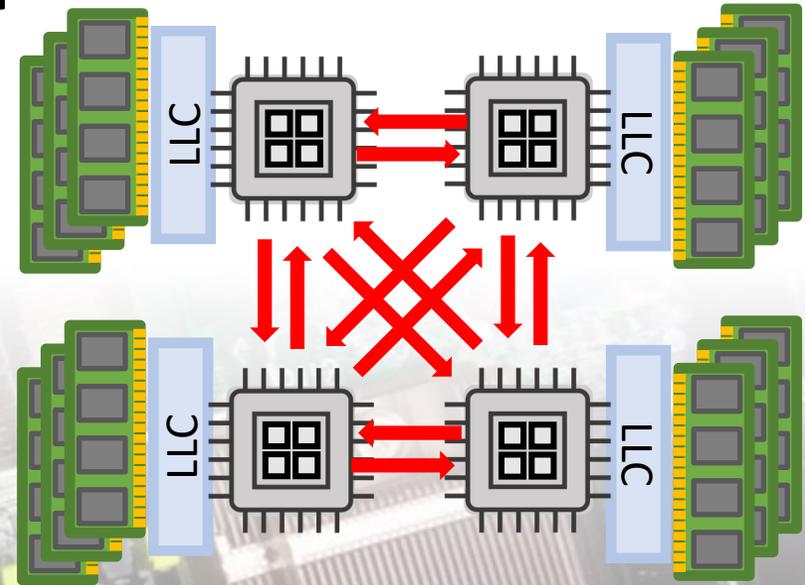
Concurrency makes it harder

- Increasing core counts
- NUMA



Concurrency makes it harder

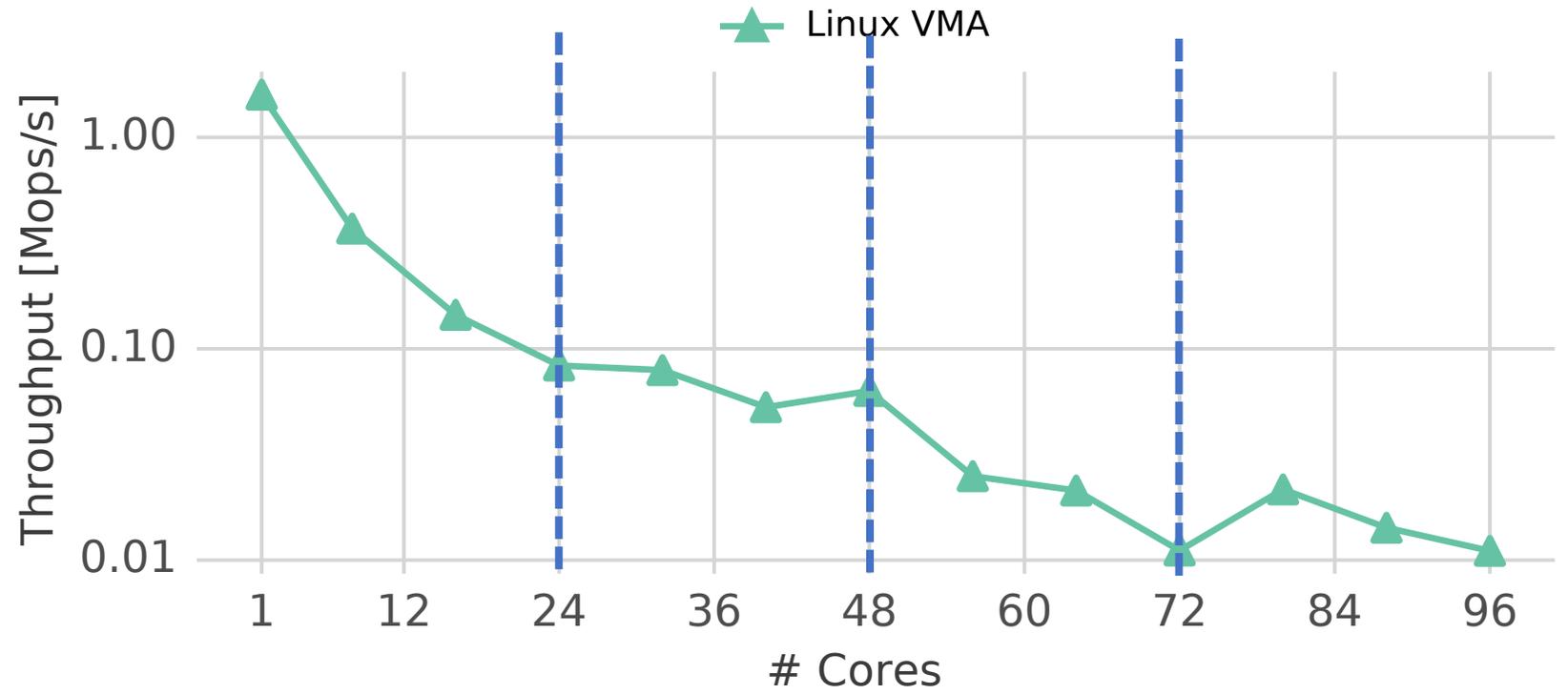
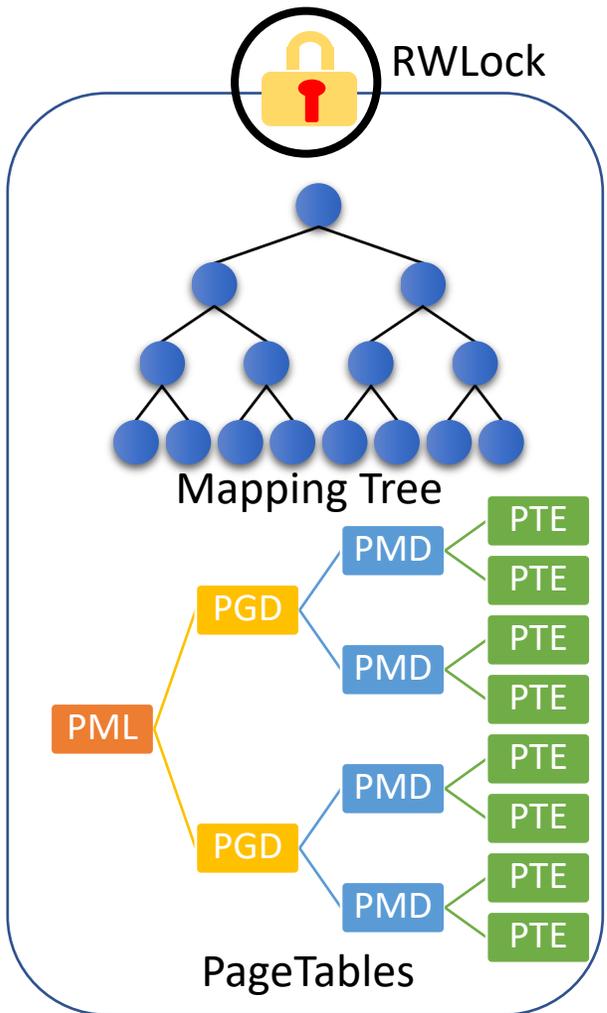
- Increasing core counts
- NUMA



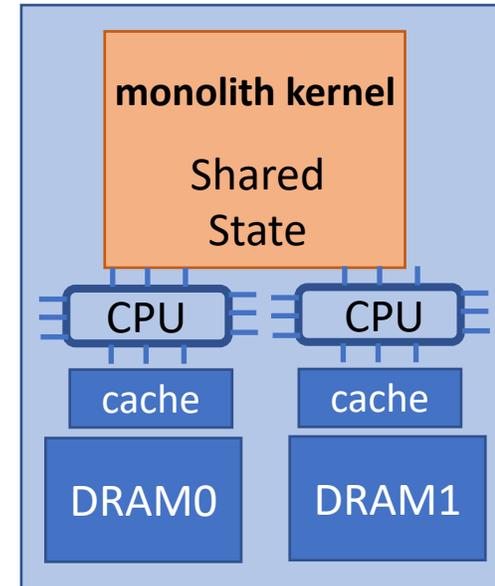
Solution: Use coarse locks around sequential data structures

Simple & black-box approach

Big locks are bad for performance



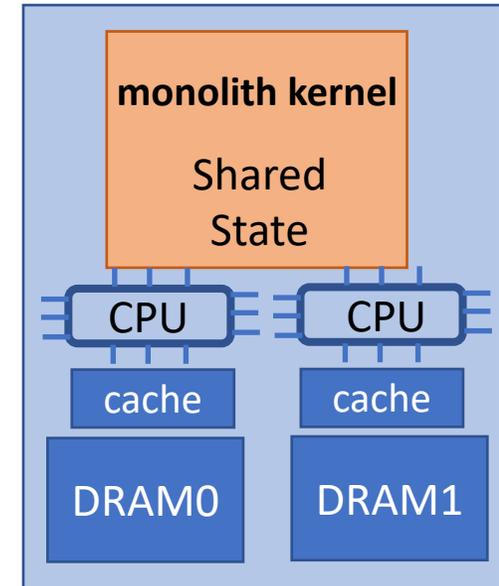
Mitigations in monolithic kernels



Mitigations in monolithic kernels

Monolithic kernels share state across cores

- Use fine-grained locking
- Lock-free data-structures
- Relativistic programming (RCU etc.)

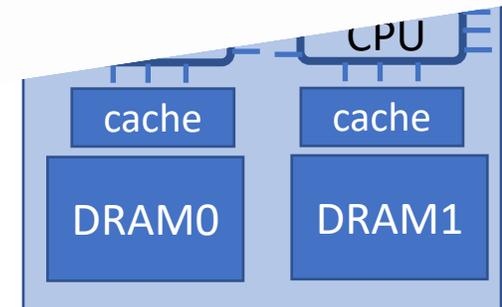


Lots of additional complexity

Mitigations in monolithic kernel

[4.18,004/235] rds: fix two RCU related problems

Message ID 20180924113104.696412478@linuxfoundation.org
State New, archived
Headers [show](#)
Series 4.18.10-stable review
Related [show](#)



Lots of additional complexity

Mitigations in monolithic kernel

^ **mm/mempool: fix two RCU related problems**
 ^ **mm/mempool: fix a data race in mempool_free()**

Message ID	1581446384-2131-1-git-send-email-cai@lca.pw
State	Accepted
Commit	abe1de4209f670ca81ba0d96f64fa9285c05a5ad
Headers	show
Series	mm/mempool: fix a data race in mempool_free()
Series	show
Related	

Lots of additional complexity

Mitigations in monolithic kernel

^ **mm/mempool: fix two RCU related problems**
 Message ID 1581446384-2131
 State Accepted
 Commit
 [v2] **ext4: fix a data race in EXT4_I(inode)->i_dsksize**

Message ID 1581085751-31793-1-git-send-email-cai@lca.pw
 State Accepted
 Commit 35df4299a6487f323b0aca120ea3f485dfee2ae3
 Headers show
 Series [v2] ext4: fix a data race in EXT4_I(inode)->i_dsksize

Lots of additional complexity

Mitigations in monolithic kernel

ext4: fix two RCU related problems

mblock: fix locking in bdev_del_partition

EXT4_I(inode)->i_disksize

Message ID: 20200901095941.2626957-1-hch@lst.de (mailing list archive)

State: New, archived

Series: [show](#)

Accepted: 35df4299a648

Series: [show](#)

[v2] ext4: fix a data race in EXT4_I(inode)->i_disksize

Lots of additional complexity

Windows

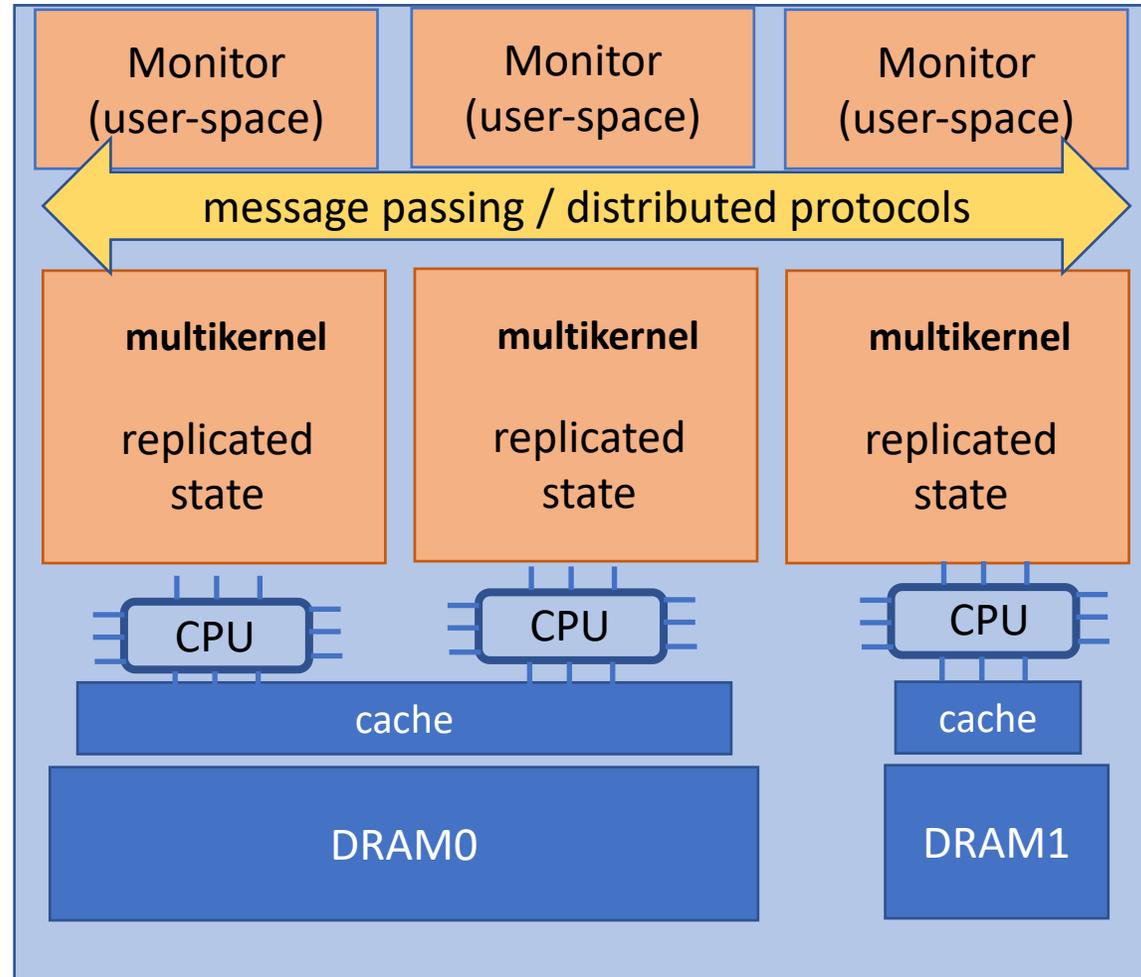
An error has occurred. We don't even know what it is. So can't fix it ,
and you have to restart your computer. By the way if you restart your ,
computer you will lose any unsaved information in all open applications.
In the other hand you don't have any other option :)

Press Enter to return to Windows (It won't work), or

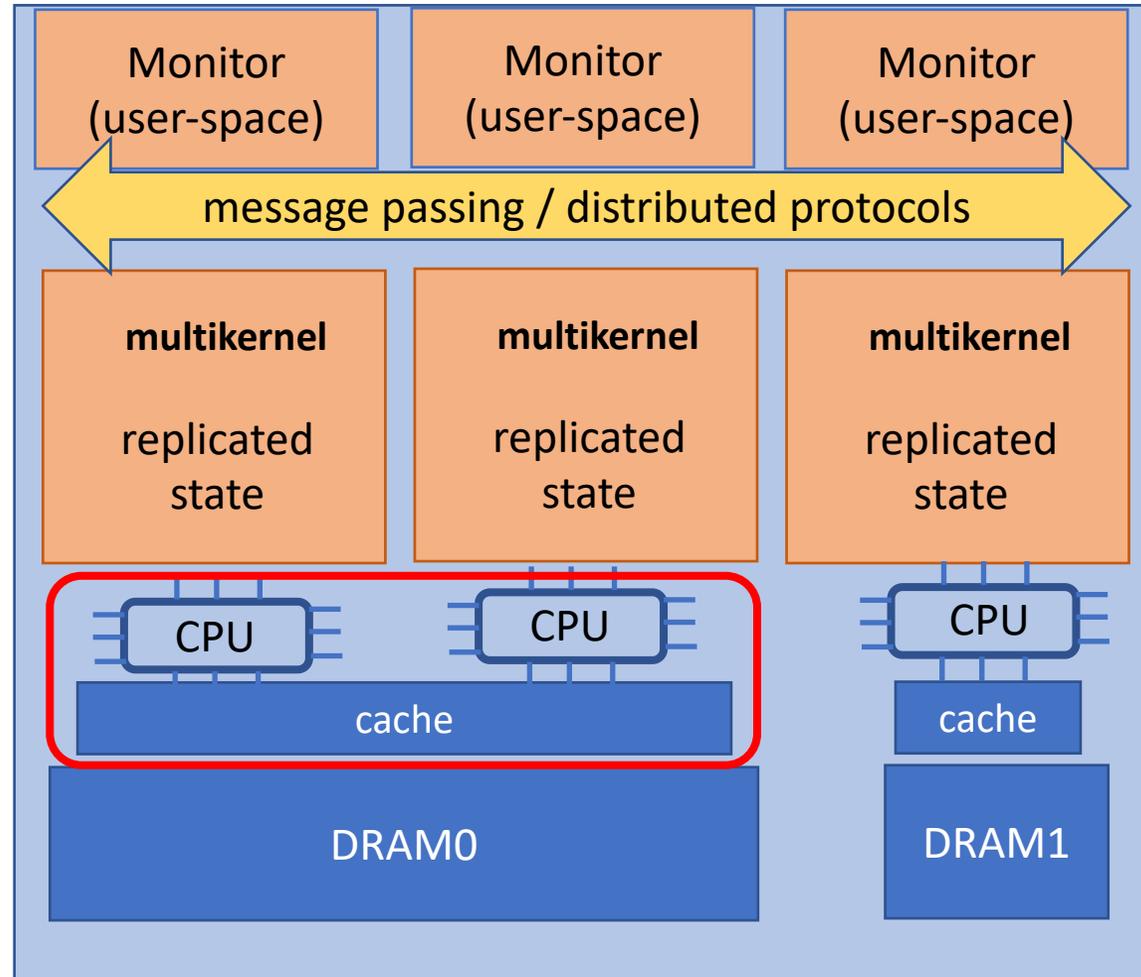
Press CTRL+ALT+DEL to restart your computer.

Error : 0E : 016F : BFF9B3D4

Multikernel: Barrelfish's approach

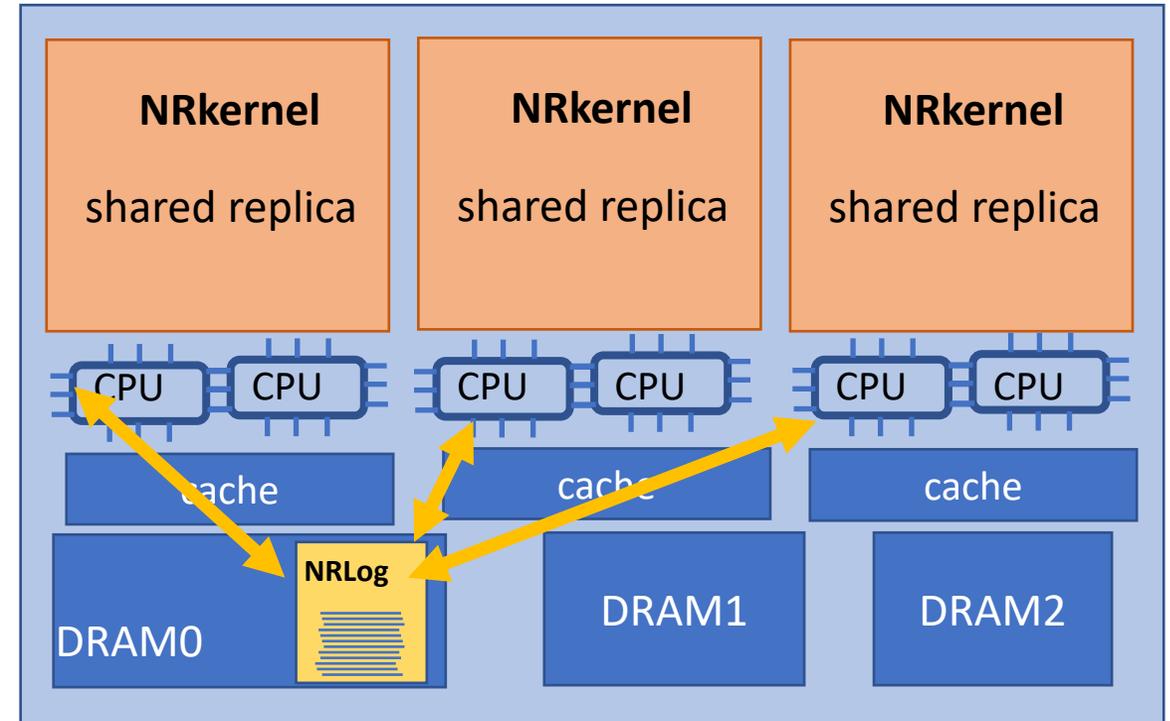


Multikernel: Barrelfish's approach



NRkernel: A log-based multikernel

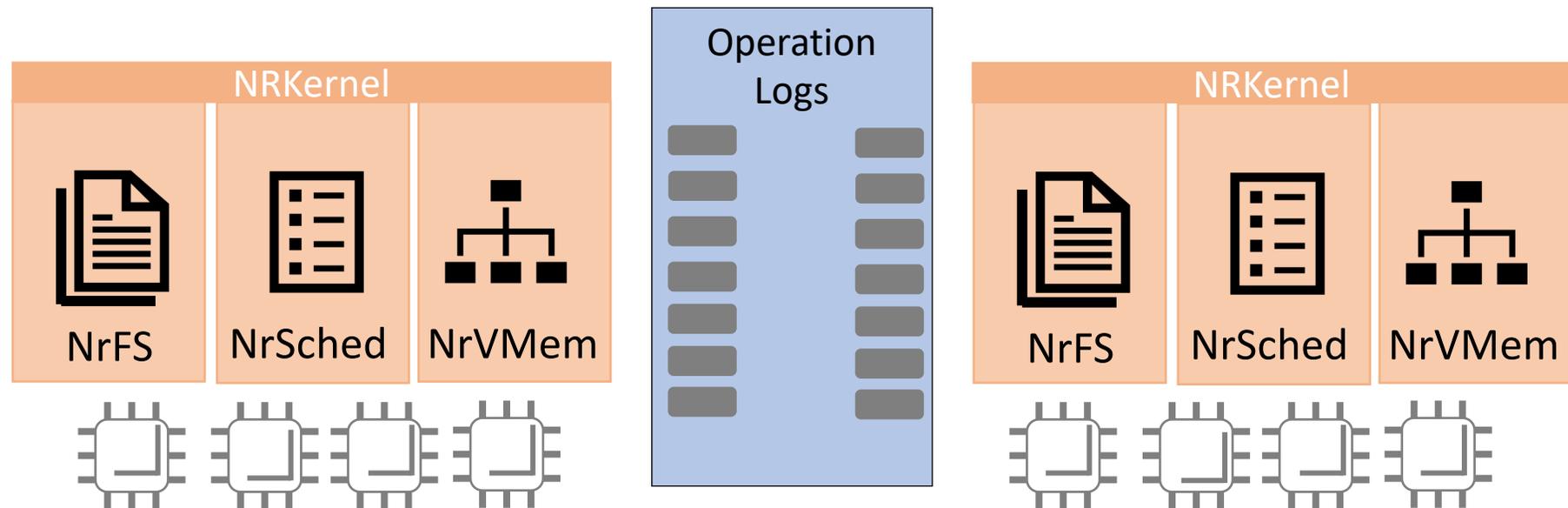
- Sequential kernel data structures
- Unified black-box approach to concurrency
- Replicates NRkernel structures
 - NUMA-local accesses
- Benefits of both message passing and shared memory



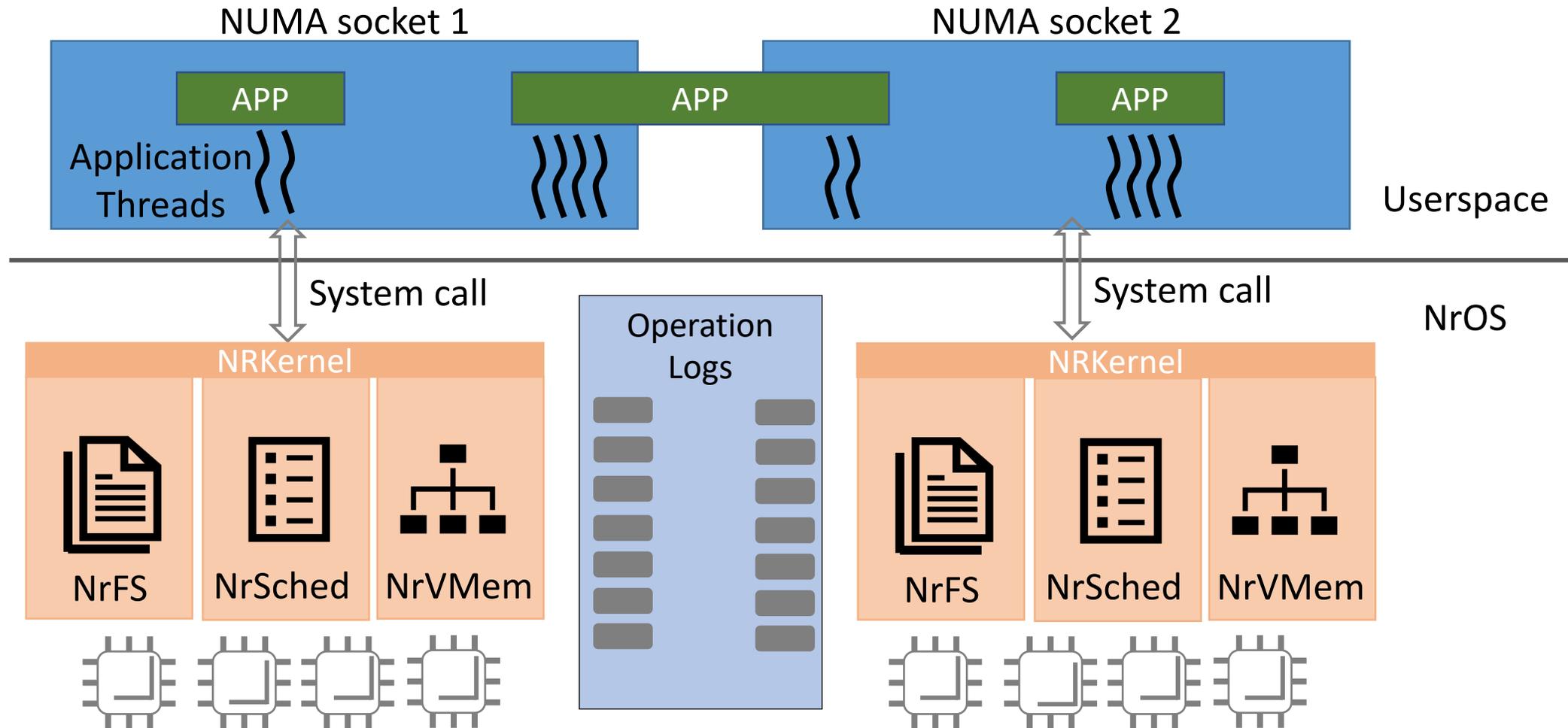
NrOS Overview

NUMA socket 1

NUMA socket 2



NrOS Overview

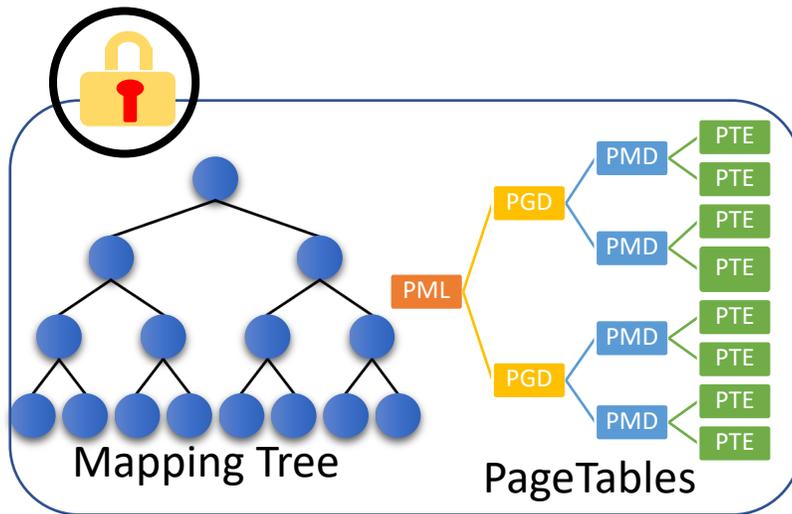


Node Replication: Log based replication

NUMA Node 1

NUMA Node 2

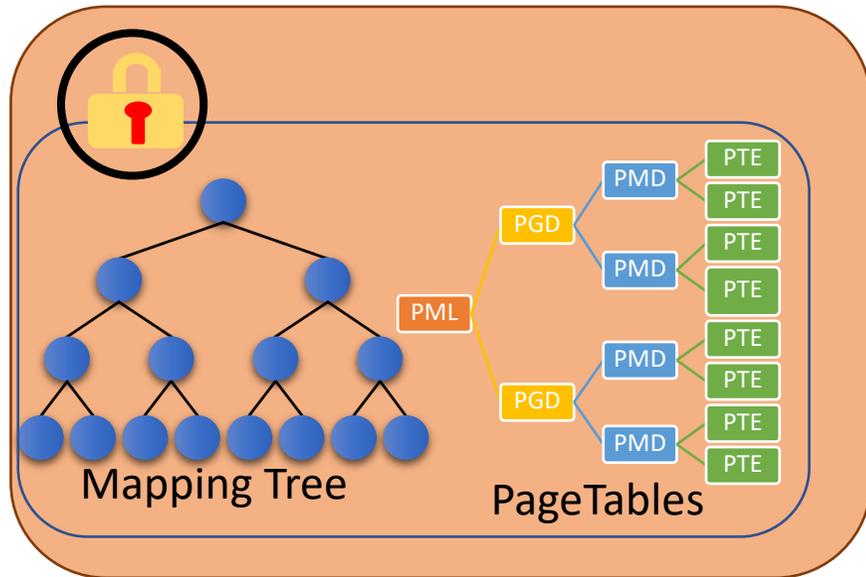
Node Replication: Log based replication



NUMA Node 1

NUMA Node 2

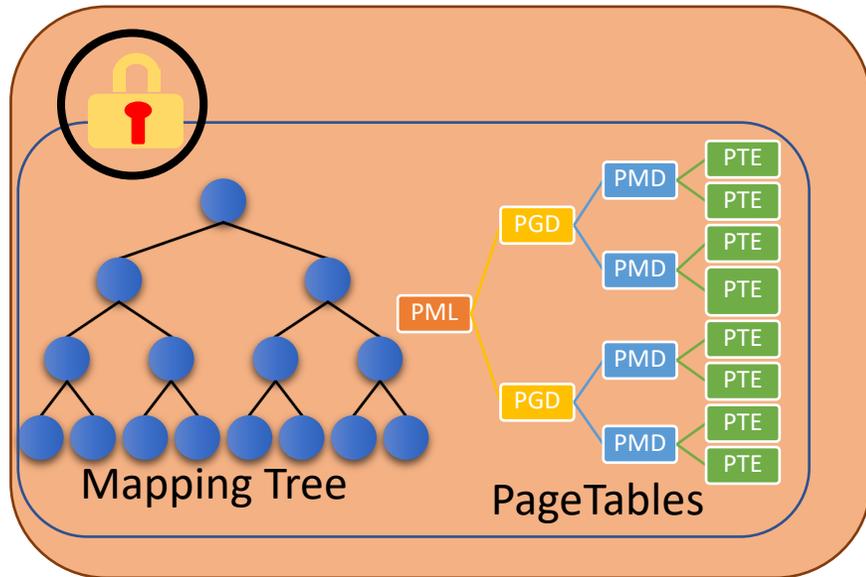
Node Replication: Log based replication



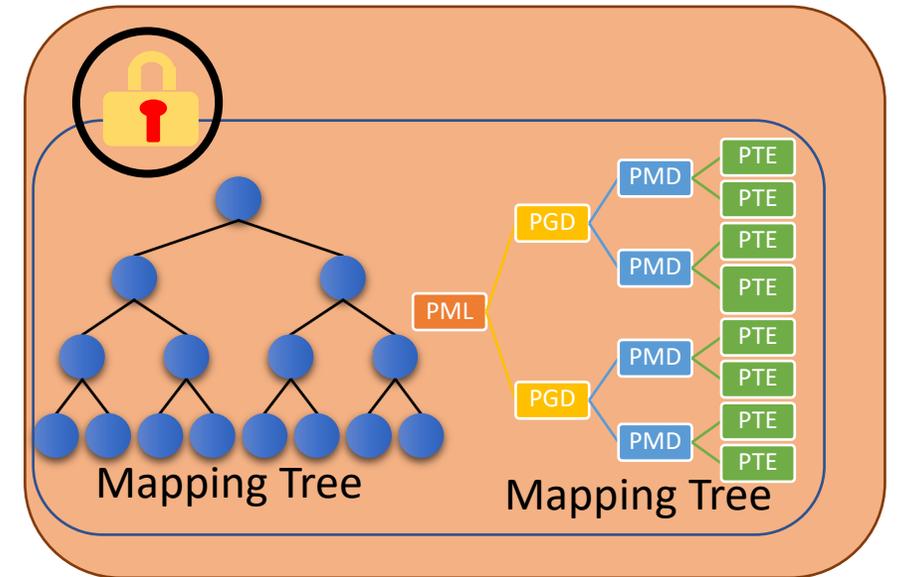
NUMA Node 1

NUMA Node 2

Node Replication: Log based replication

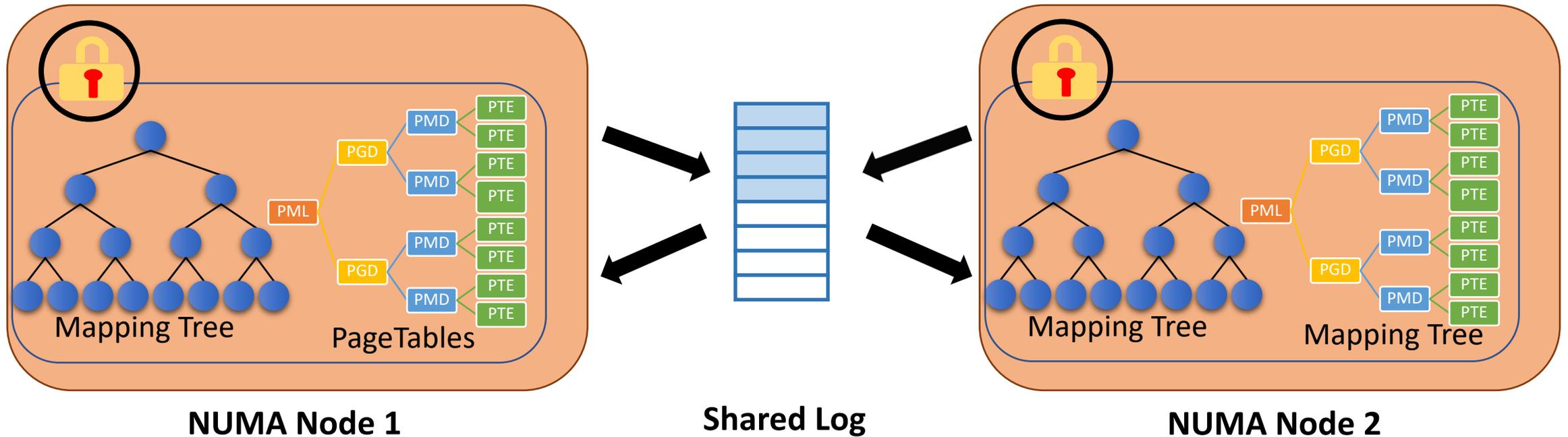


NUMA Node 1

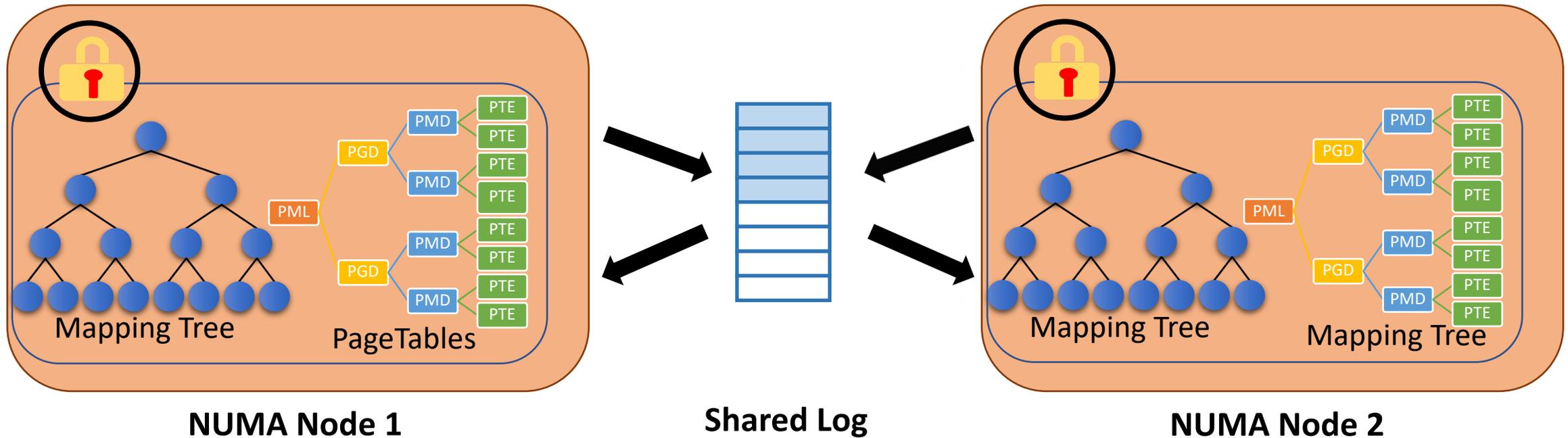


NUMA Node 2

Node Replication: Log based replication

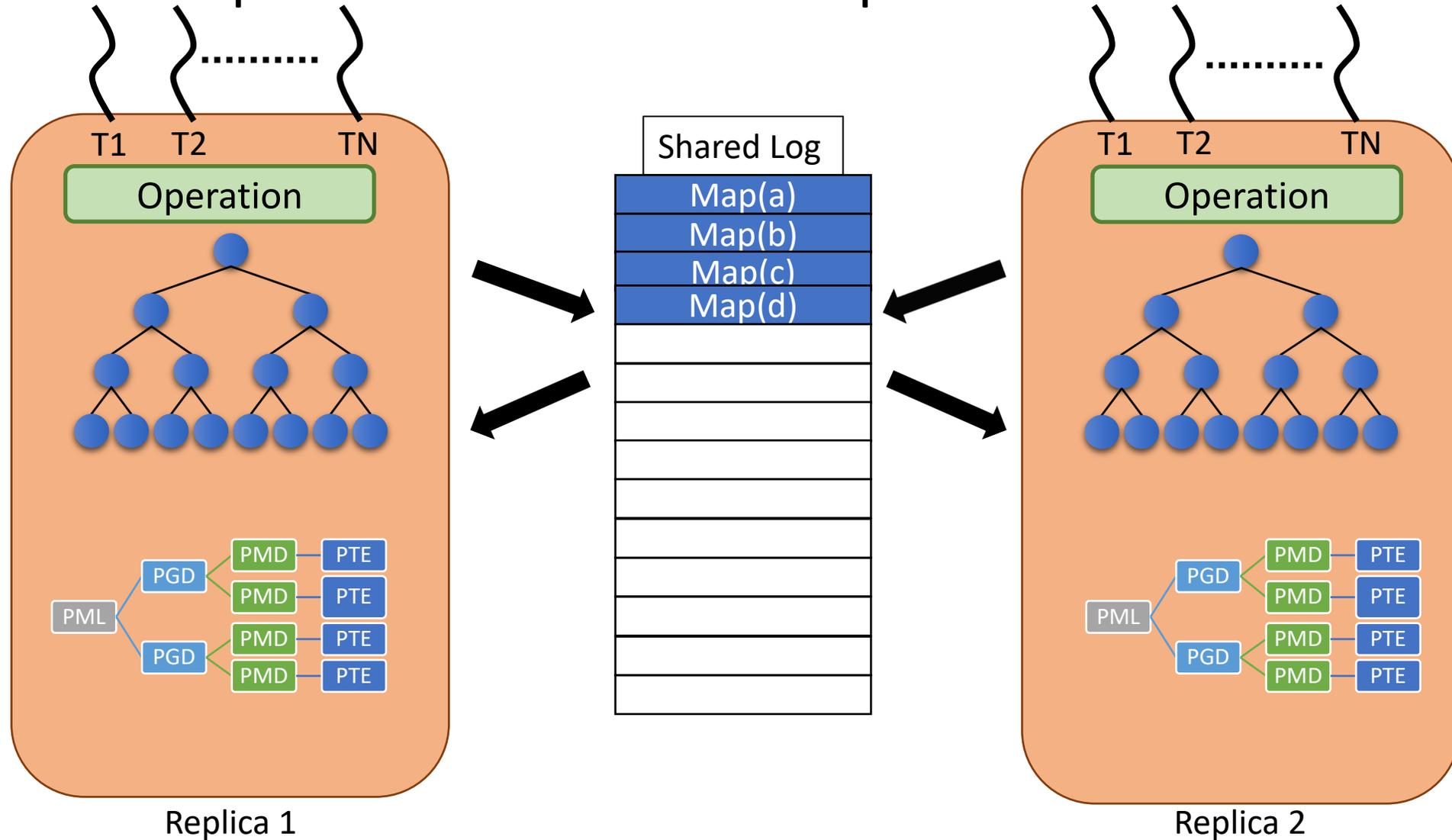


Node Replication: Log based replication

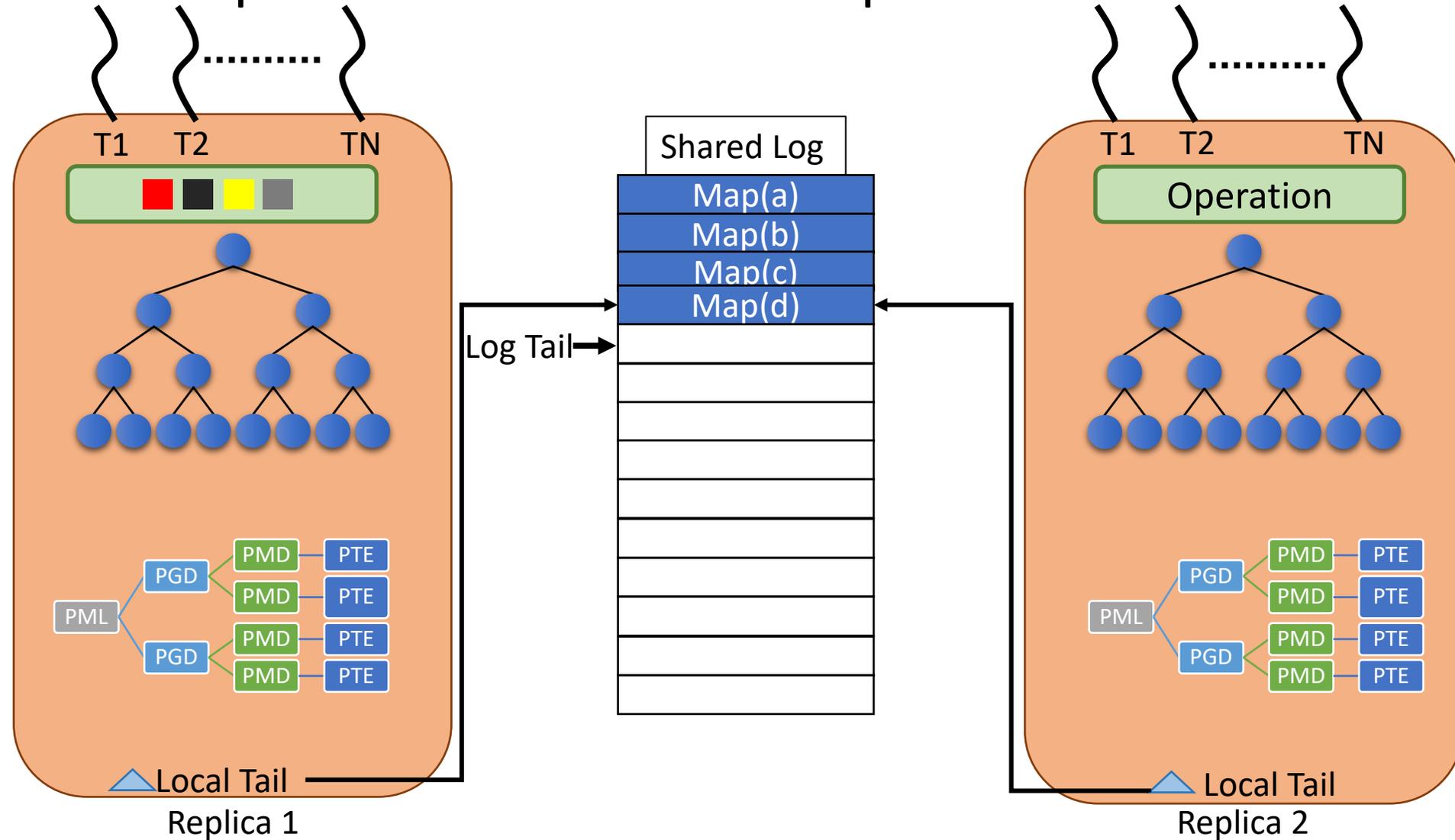


- Write operations replicate through the shared log
- Read operations access the local replica

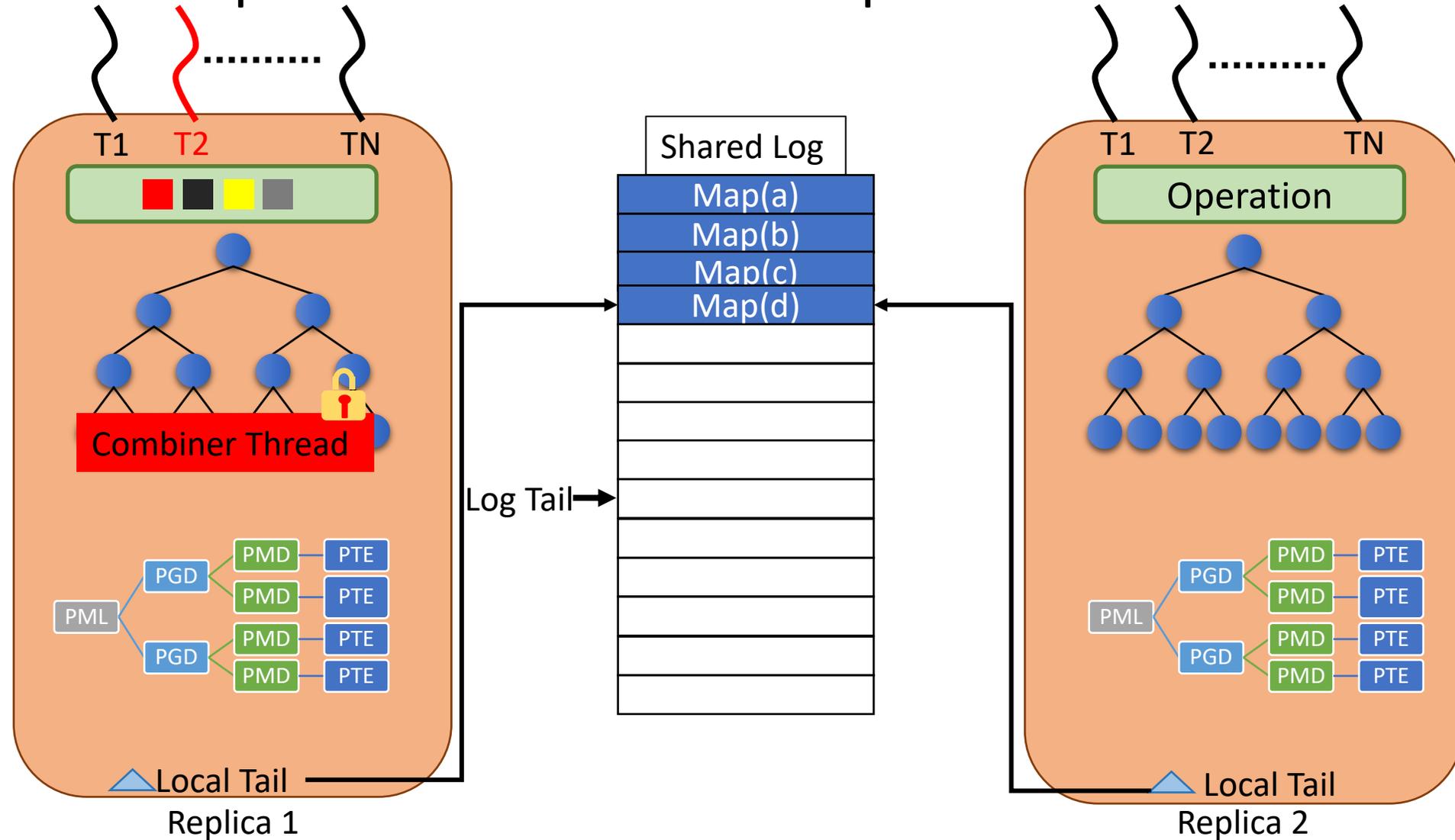
Node Replication: Write Operation



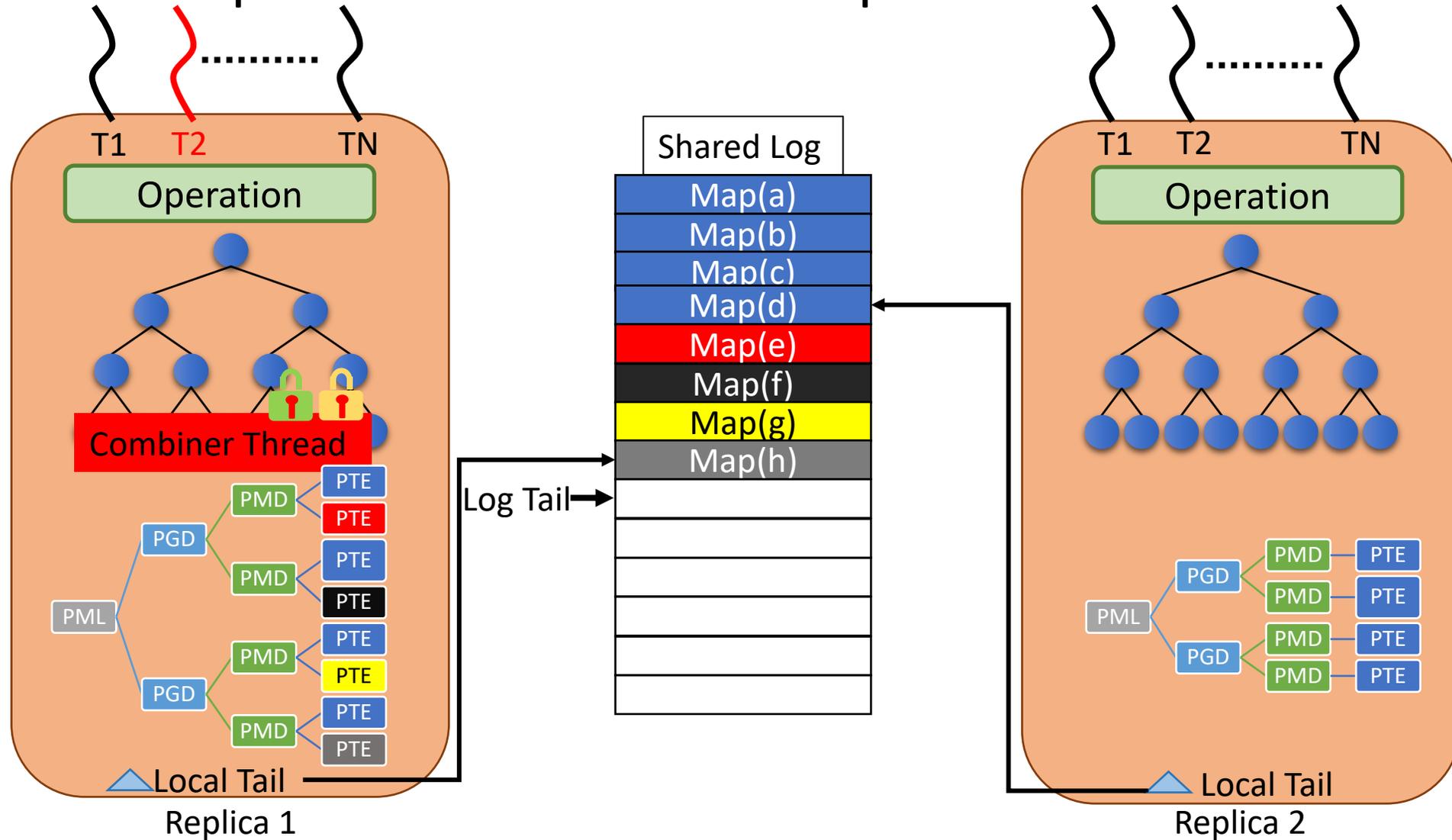
Node Replication: Write Operation



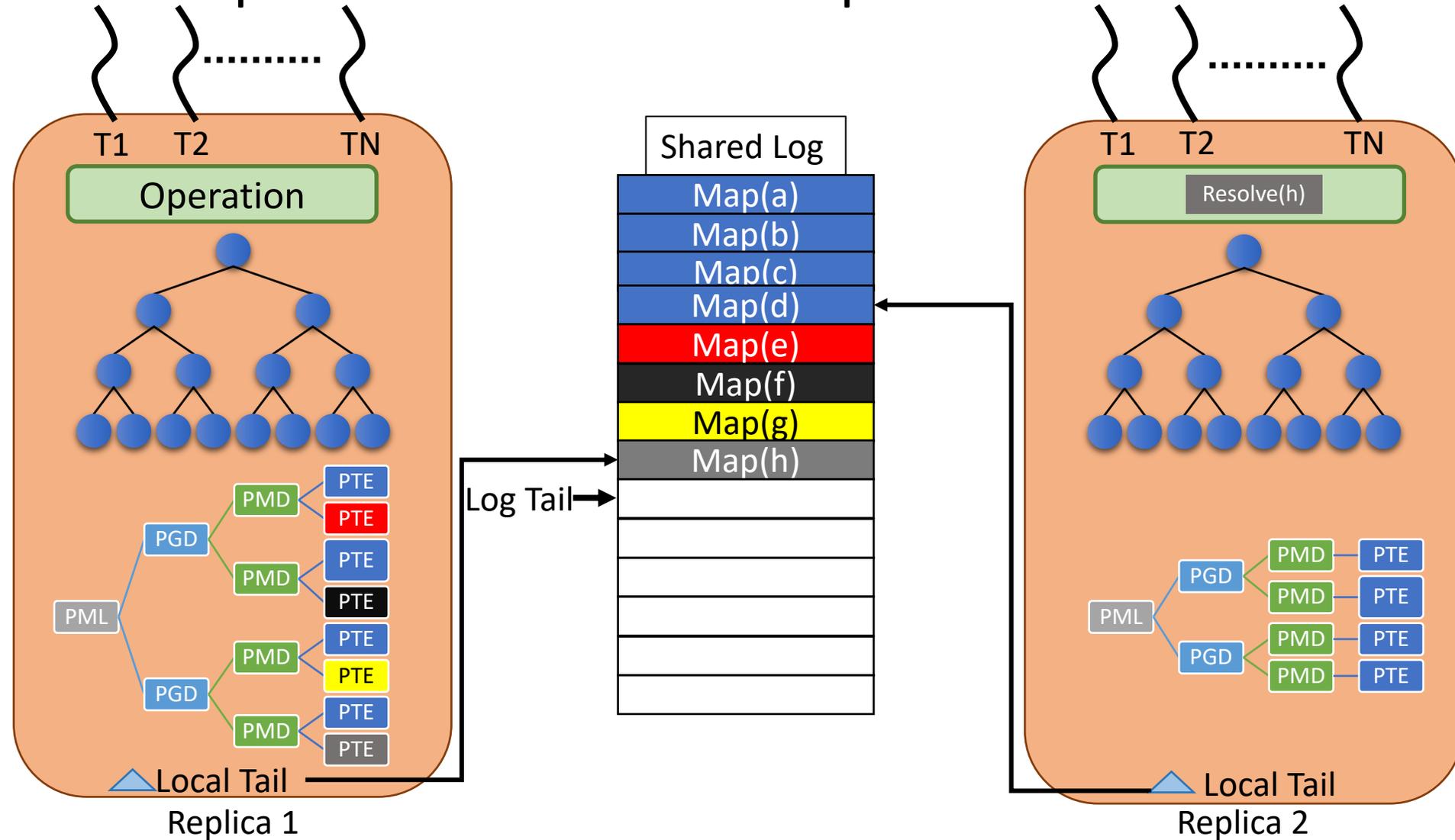
Node Replication: Write Operation



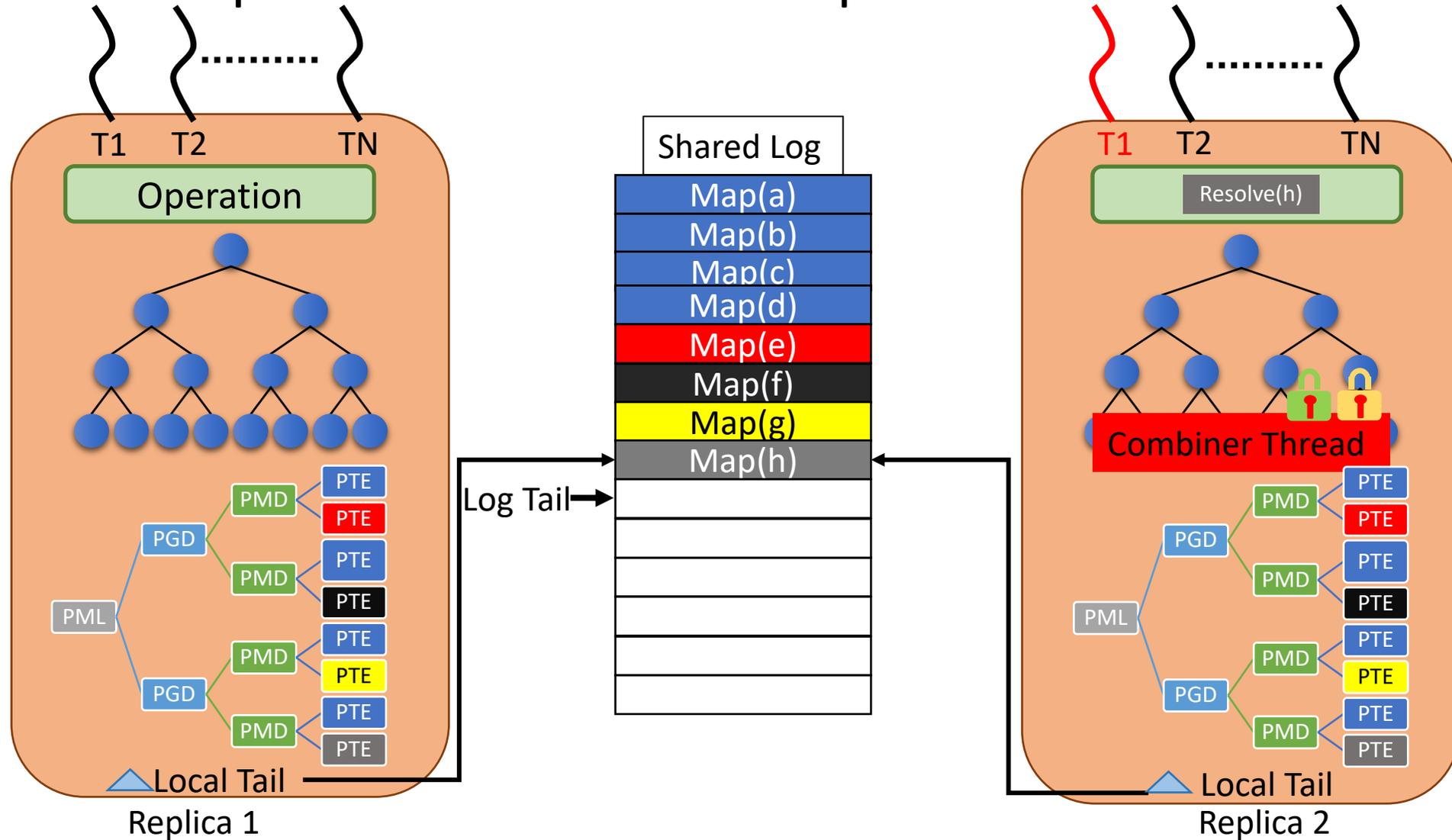
Node Replication: Write Operation



Node Replication: Read Operation



Node Replication: Read Operation



Node Replication: Read Operation

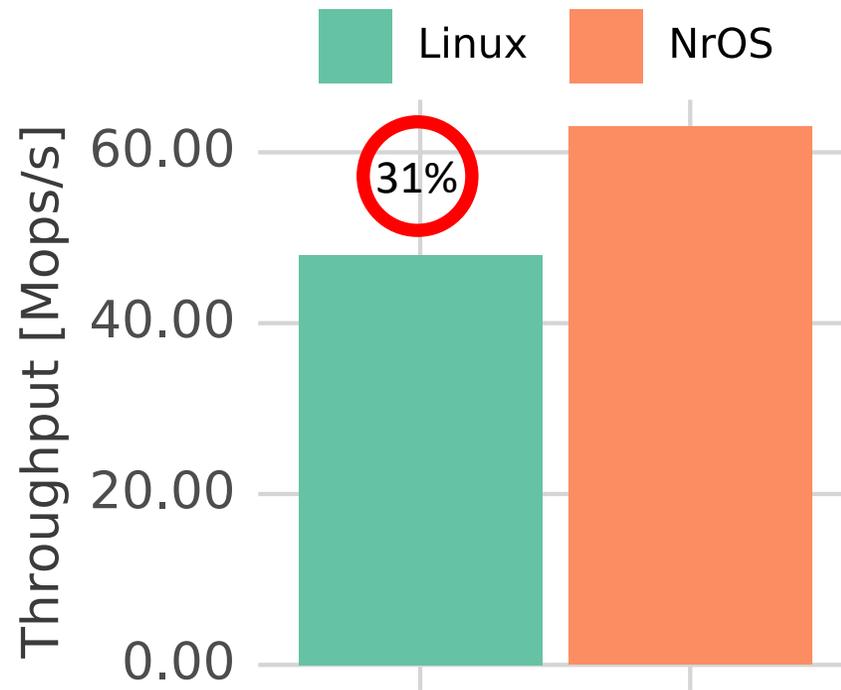


Benefits

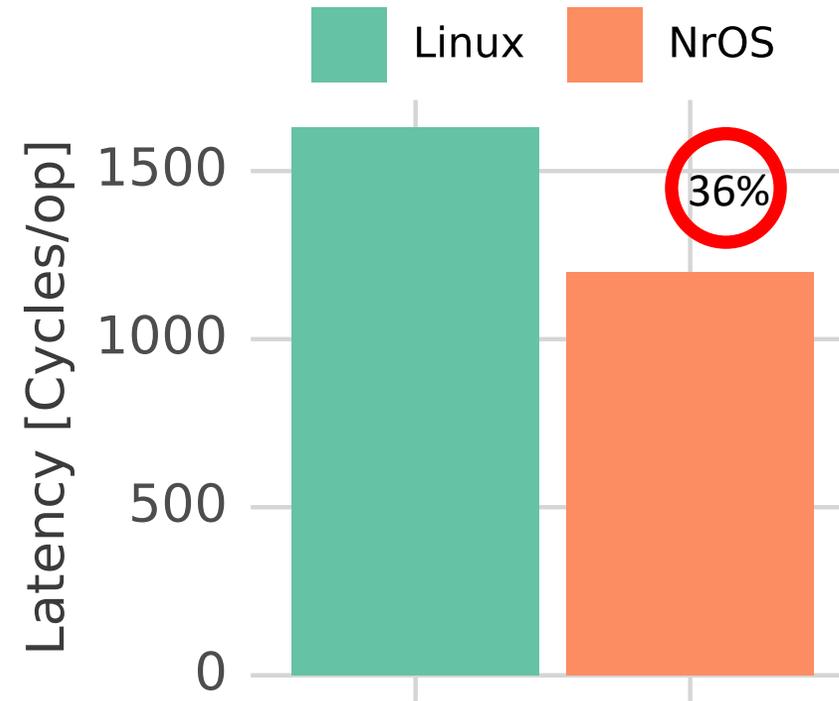
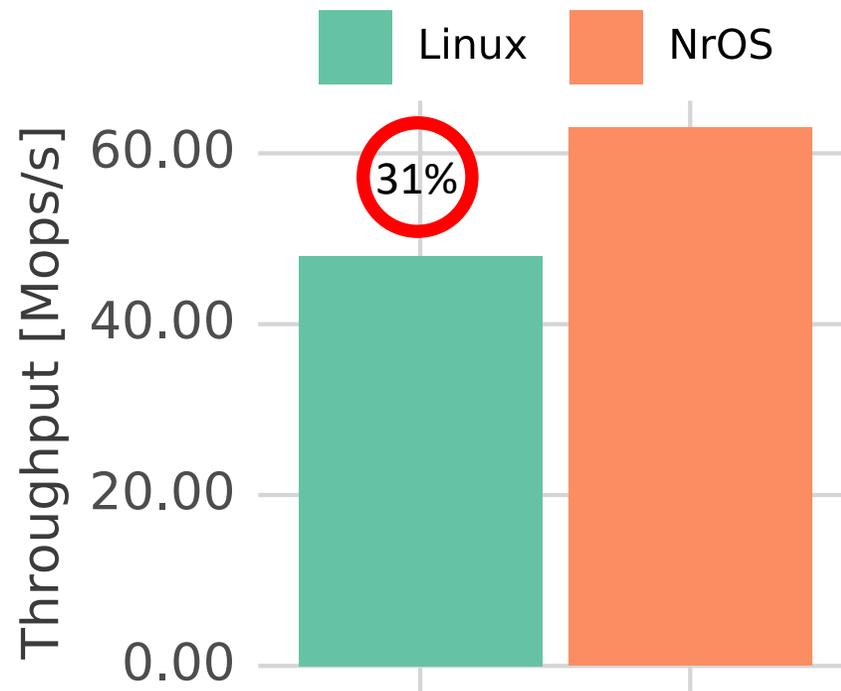
- Safe by construction: black-box technique
- NUMA-local concurrent reads
- Avoids contention with batched writes per replica (flat-combining)



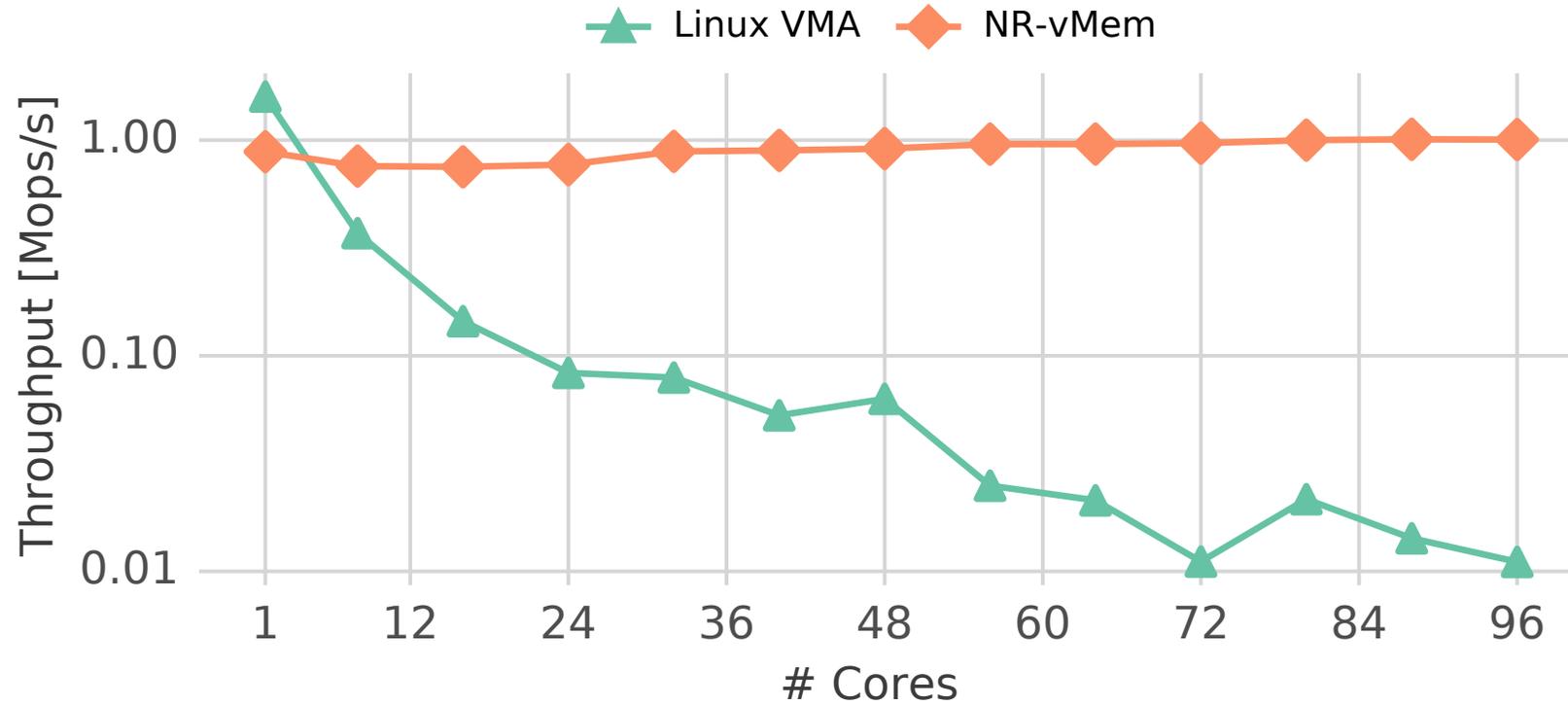
Memcached: Page-table replication benefits



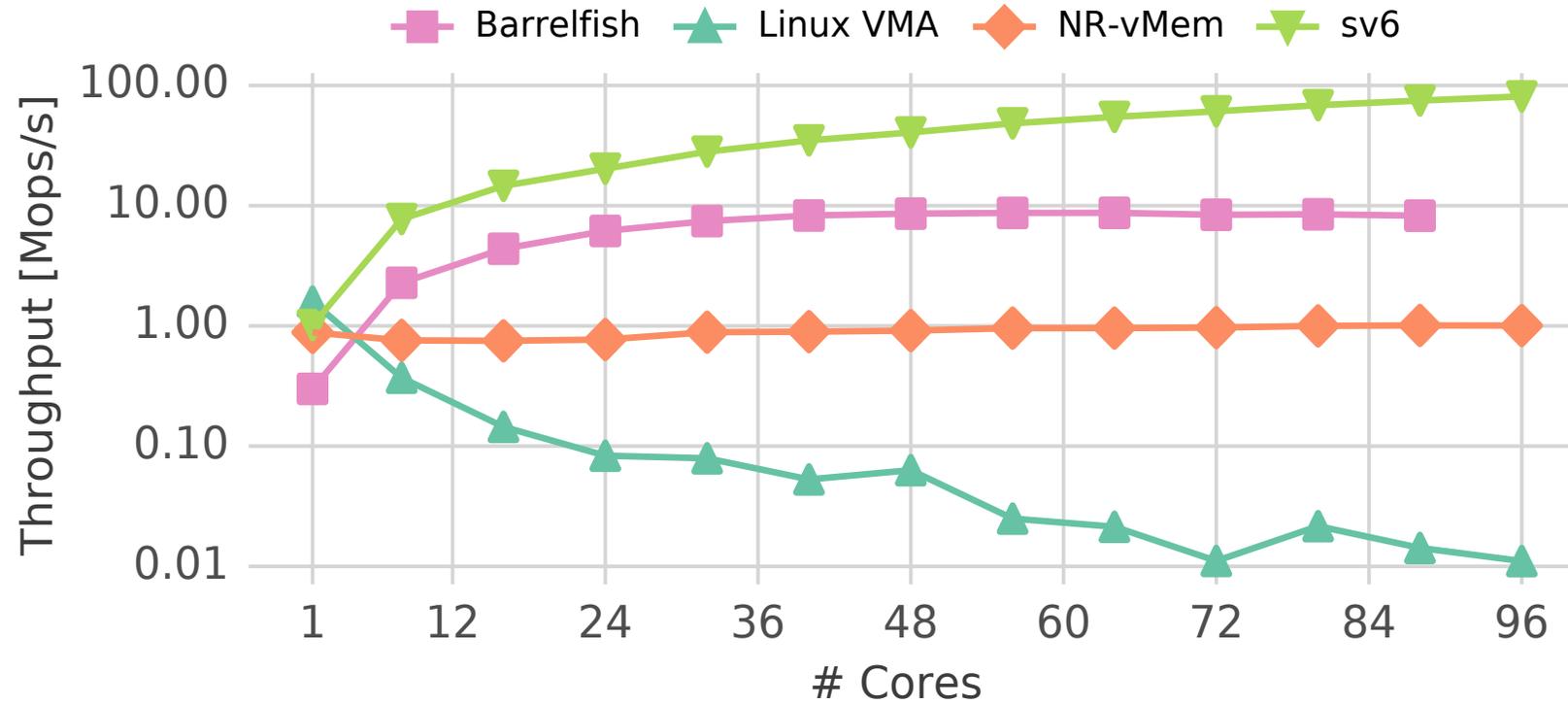
Memcached: Page-table replication benefits



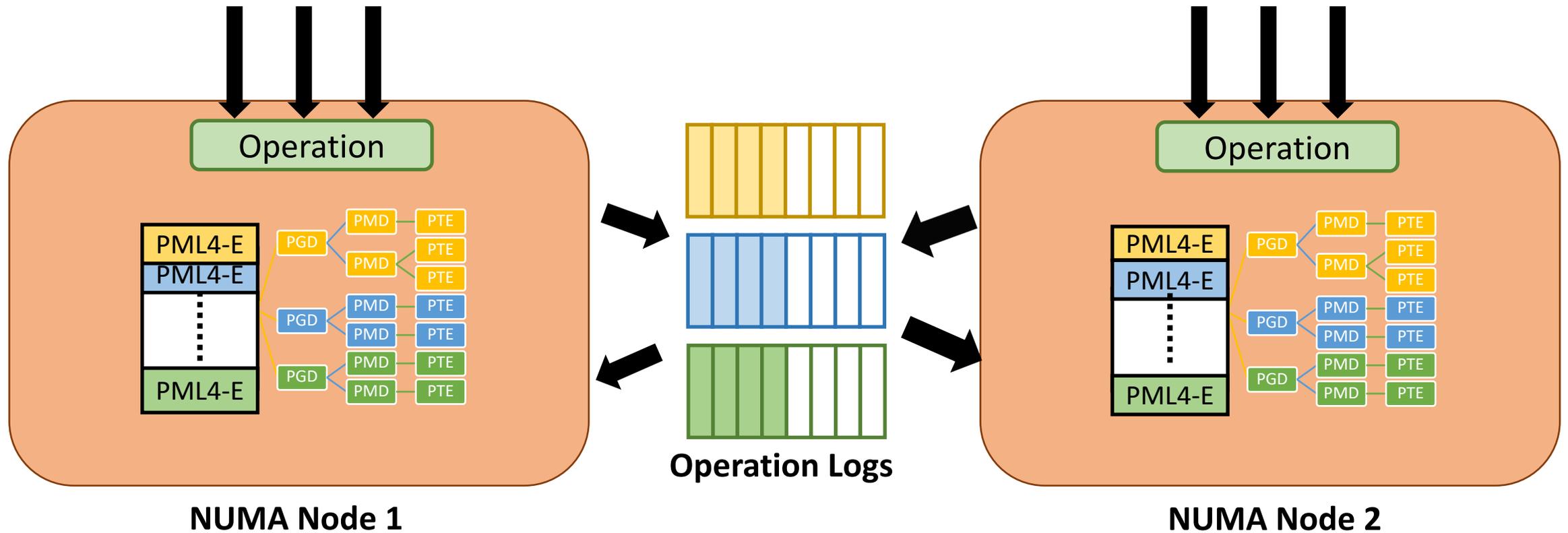
NR-VMem: Combining helps



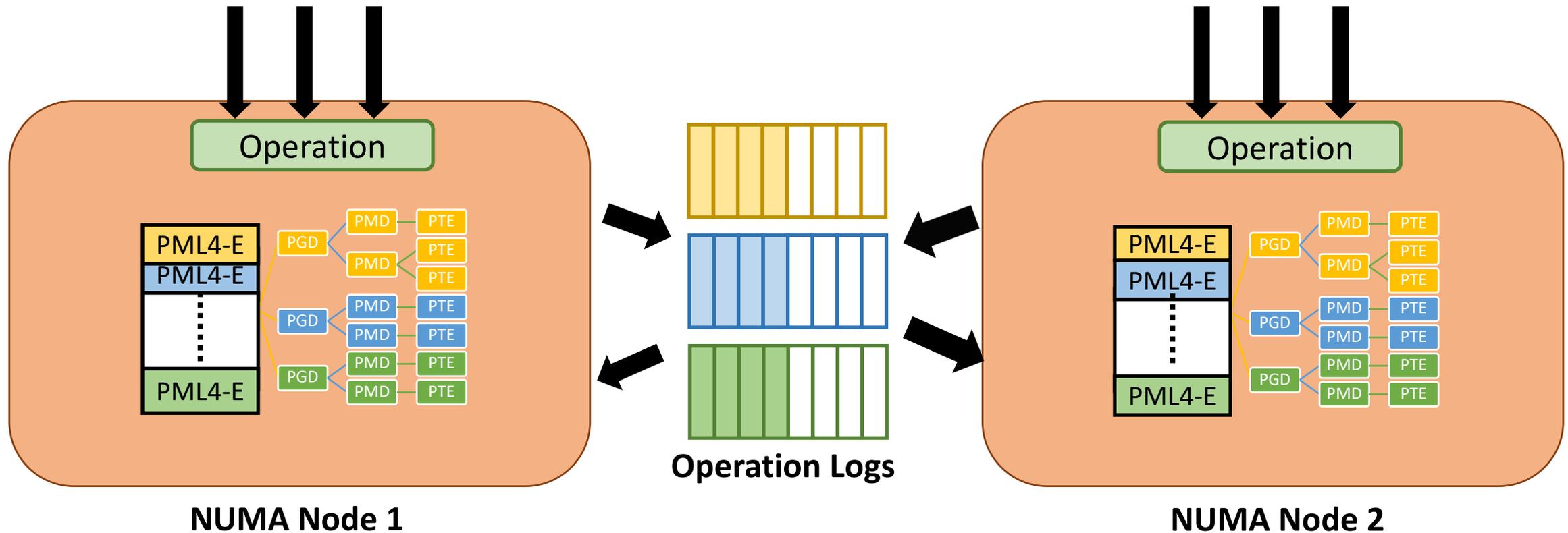
NR-VMem: Combining helps



Concurrent Node Replication (CNR)

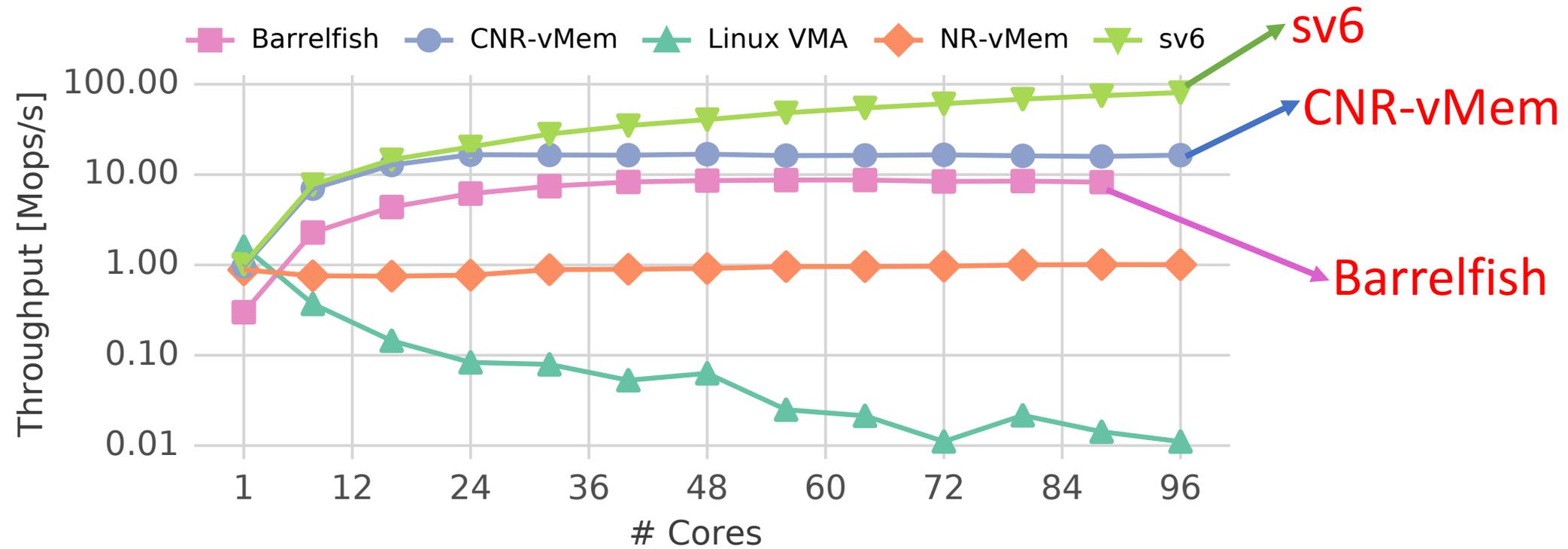


Concurrent Node Replication (CNR)

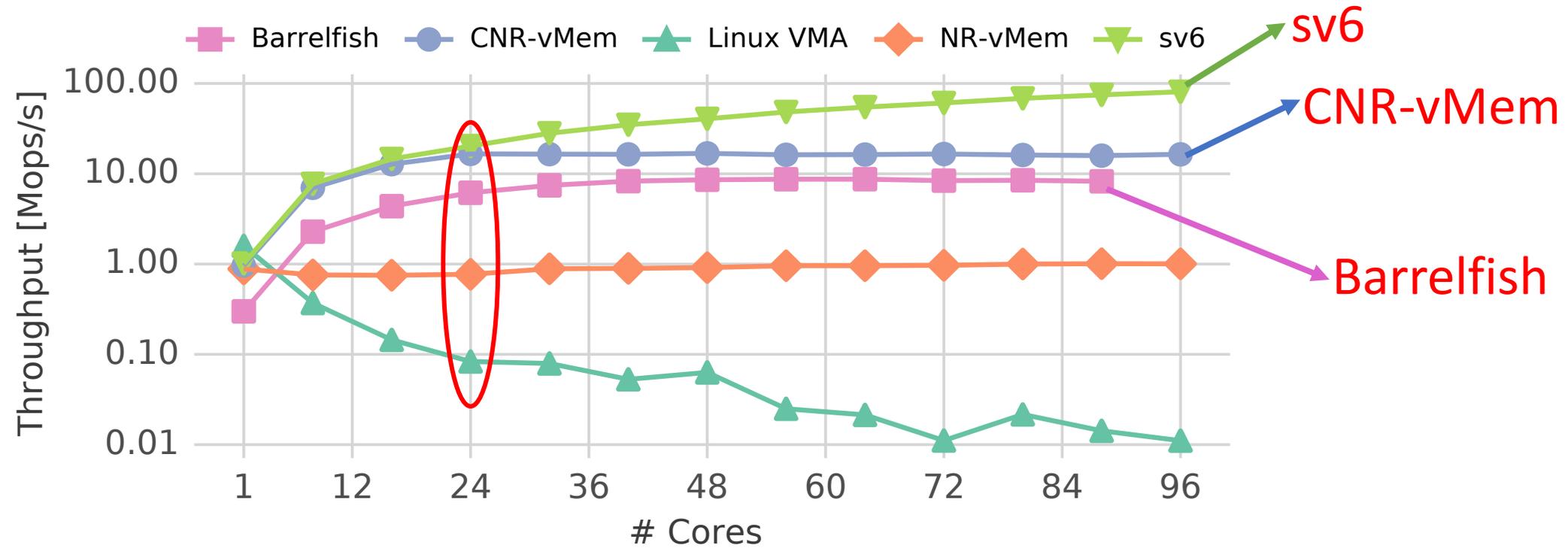


- Commutative operations can map to different logs
- Conflicting operations map to the same log

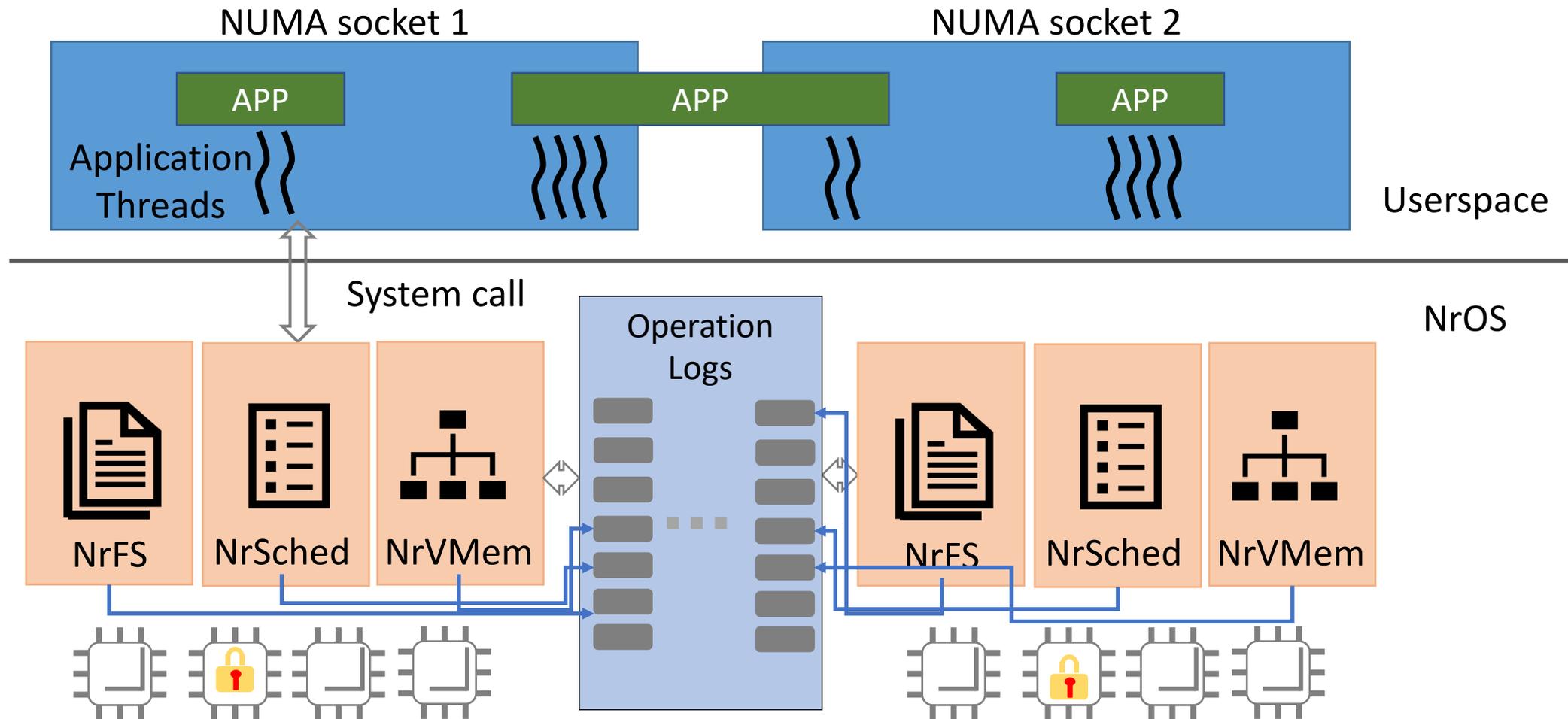
CNR-VMem: Scalable, NUMA-local PageTables



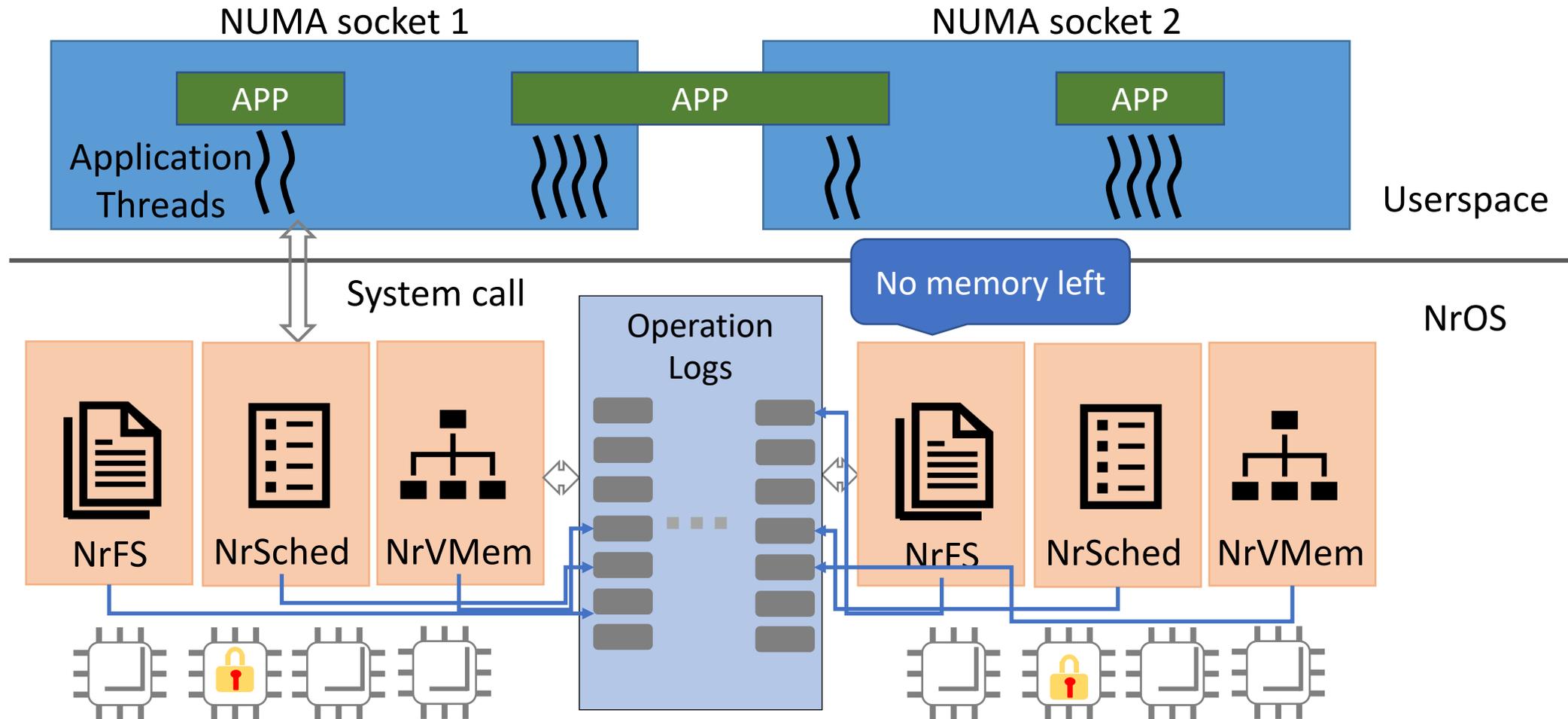
CNR-VMem: Scalable, NUMA-local PageTables



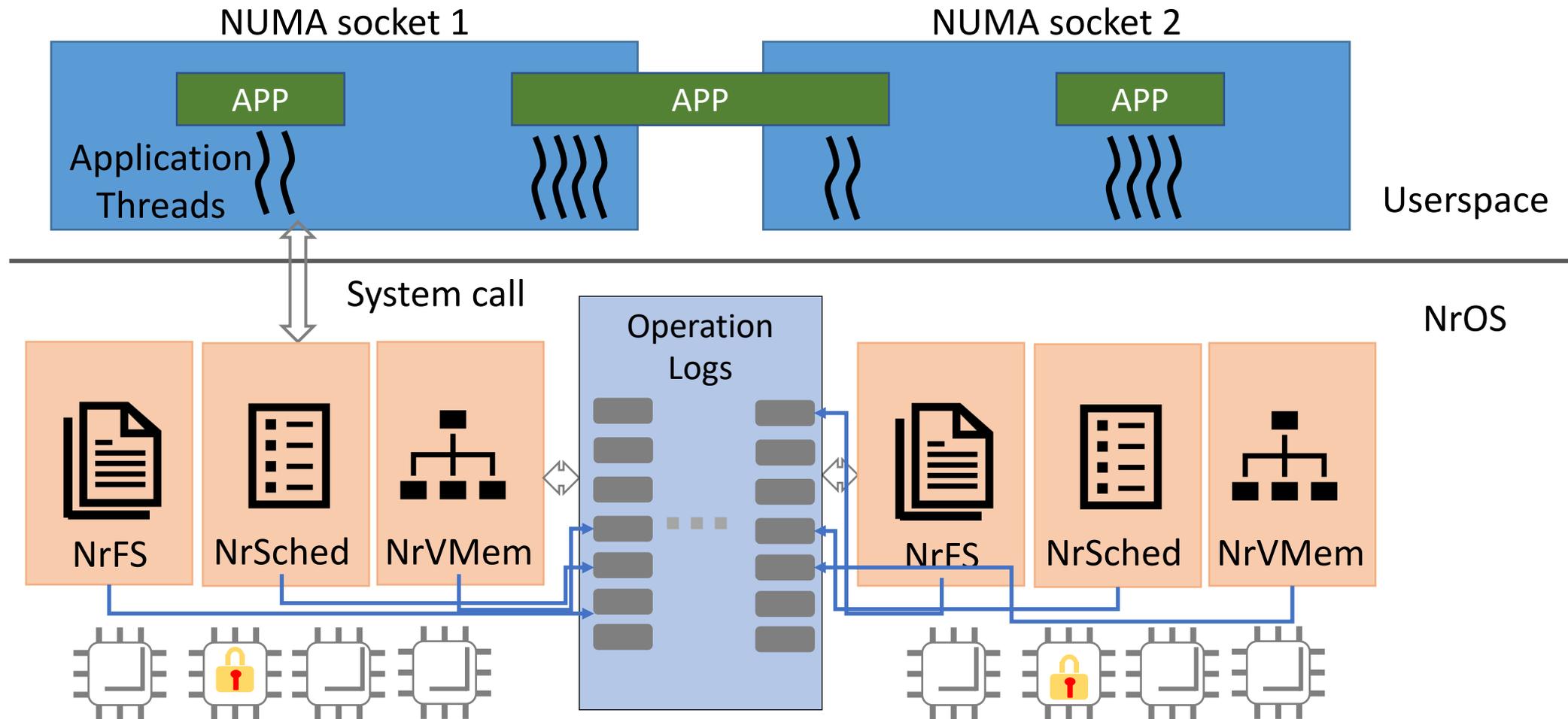
NrOS Summary



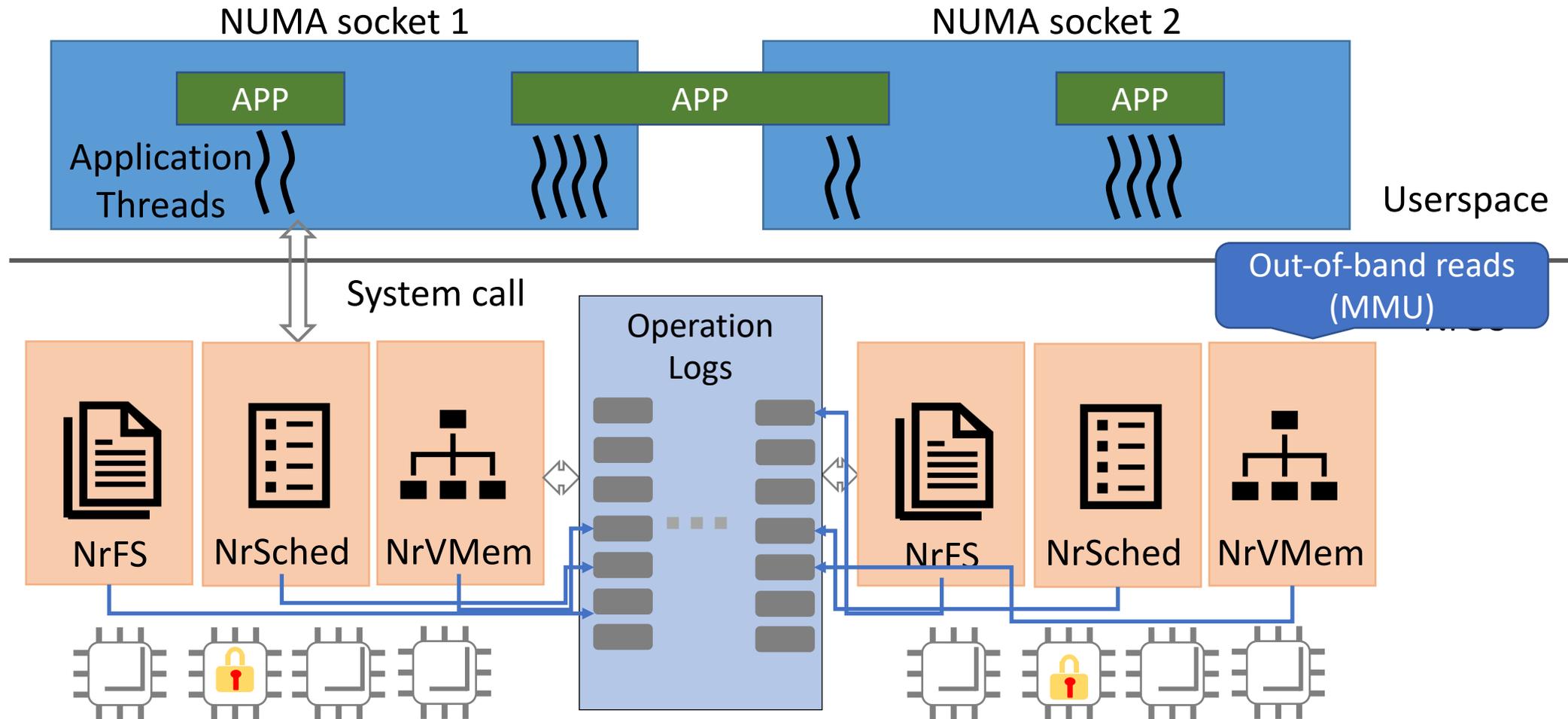
NrOS Summary



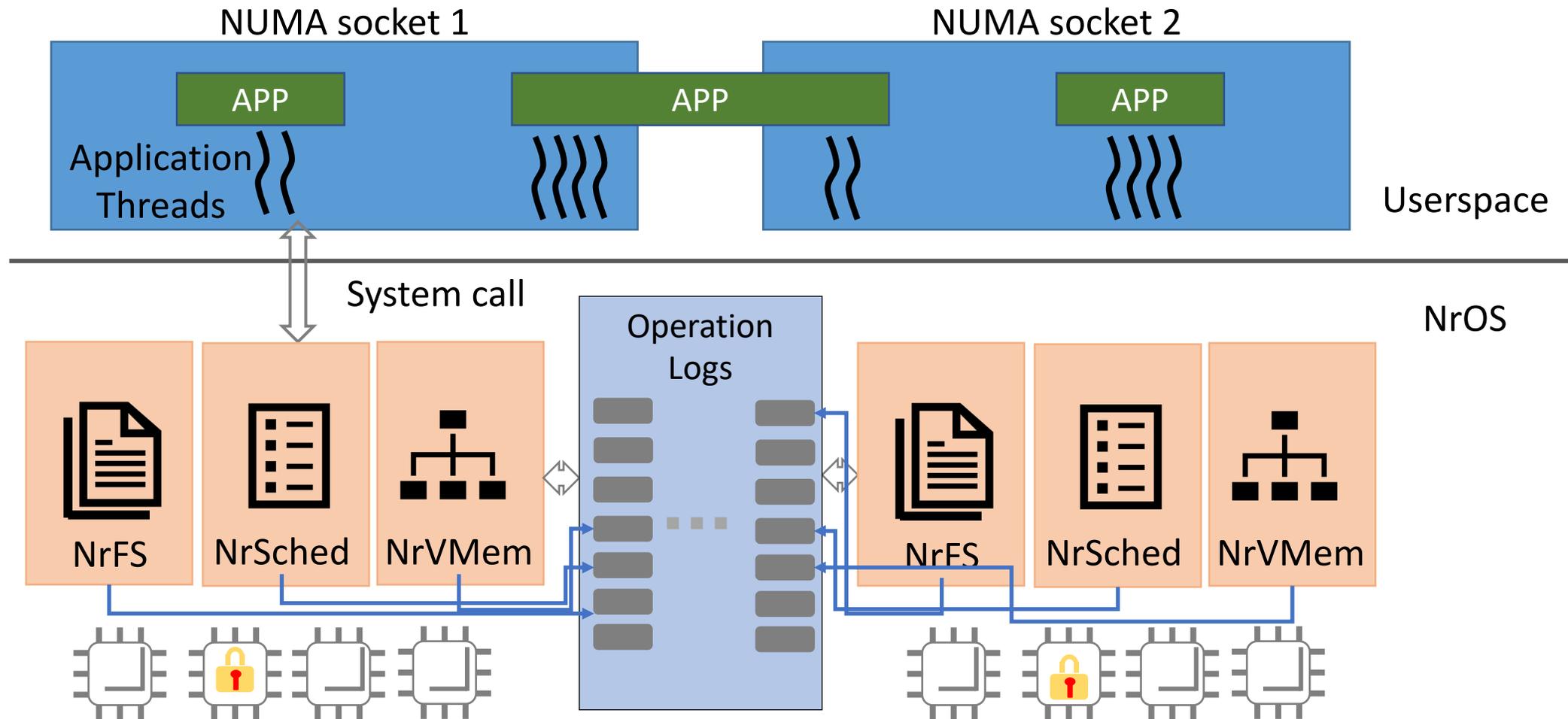
NrOS Summary



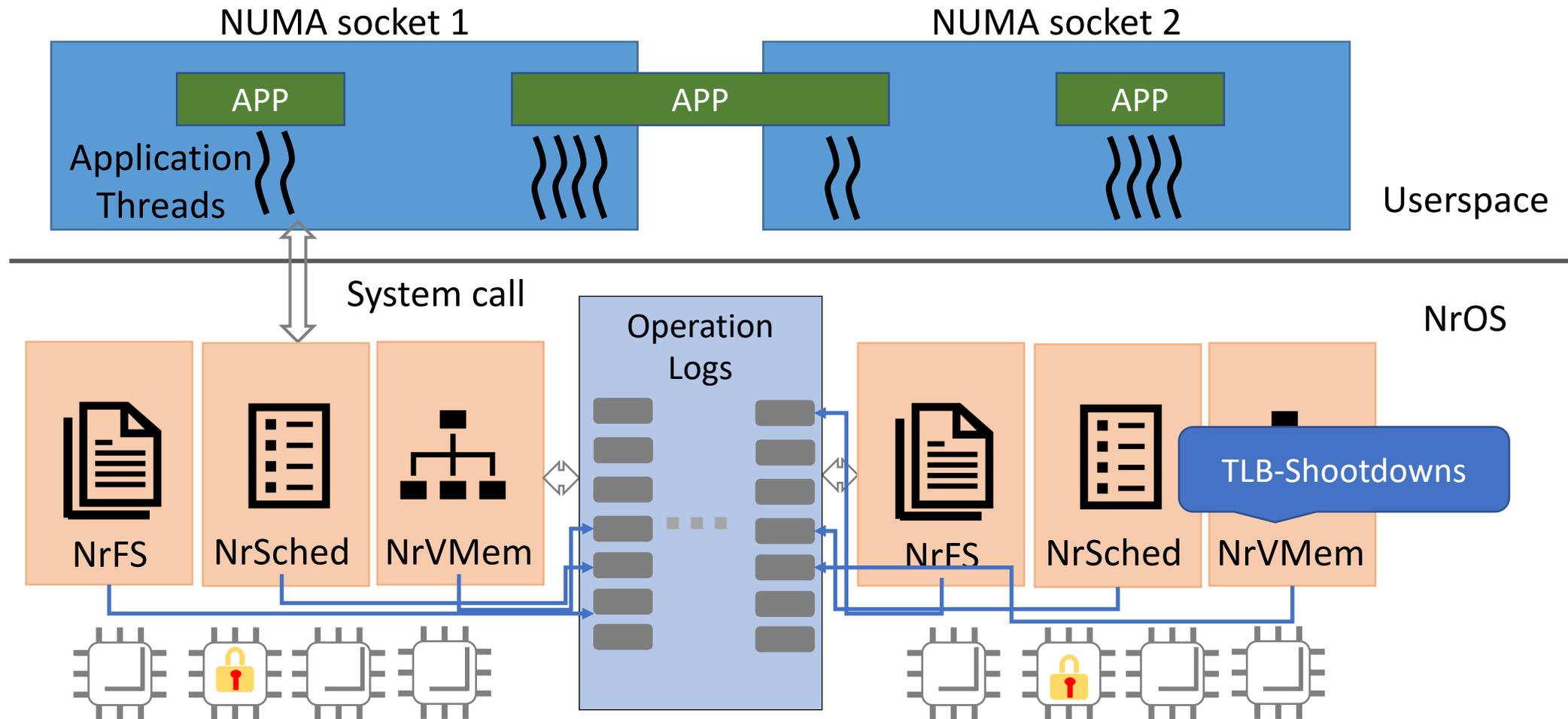
NrOS Summary



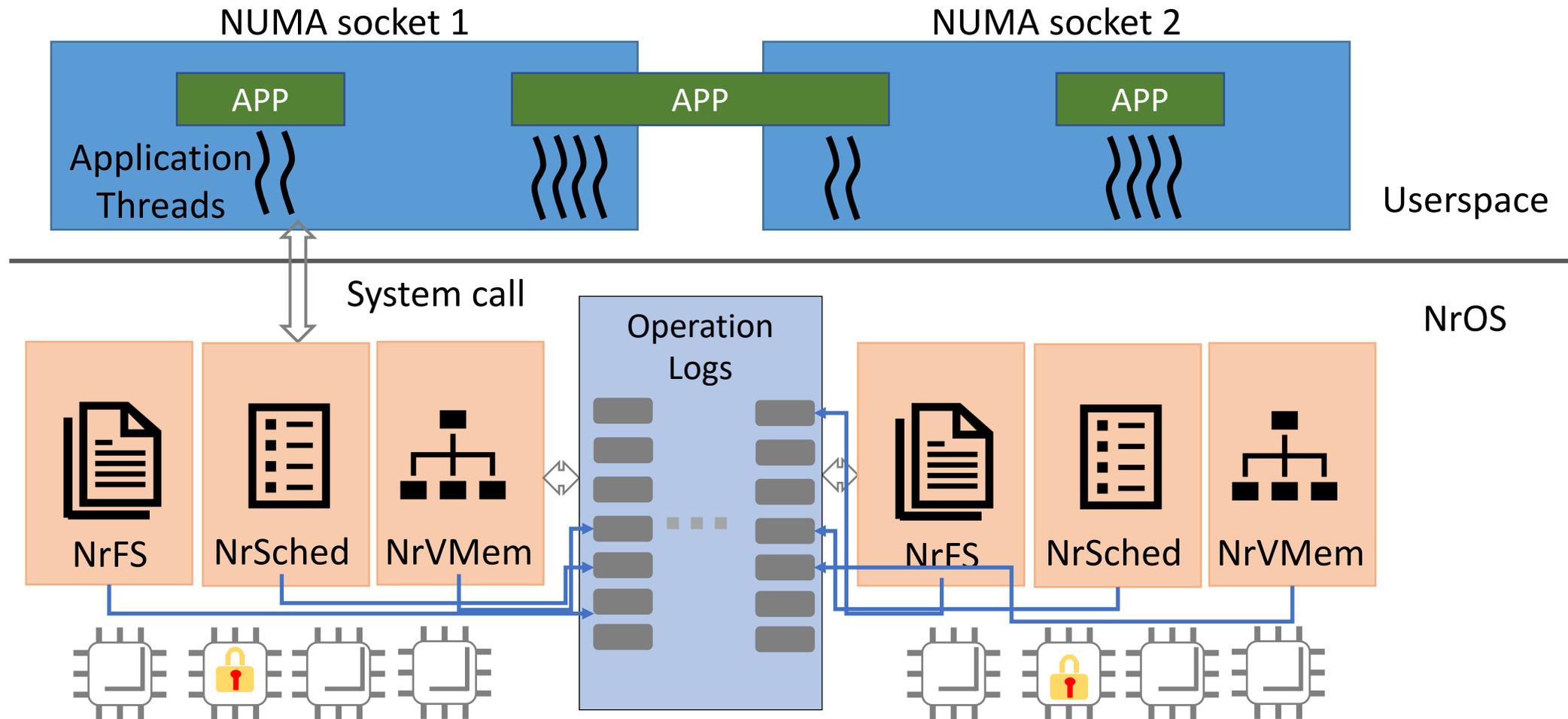
NrOS Summary



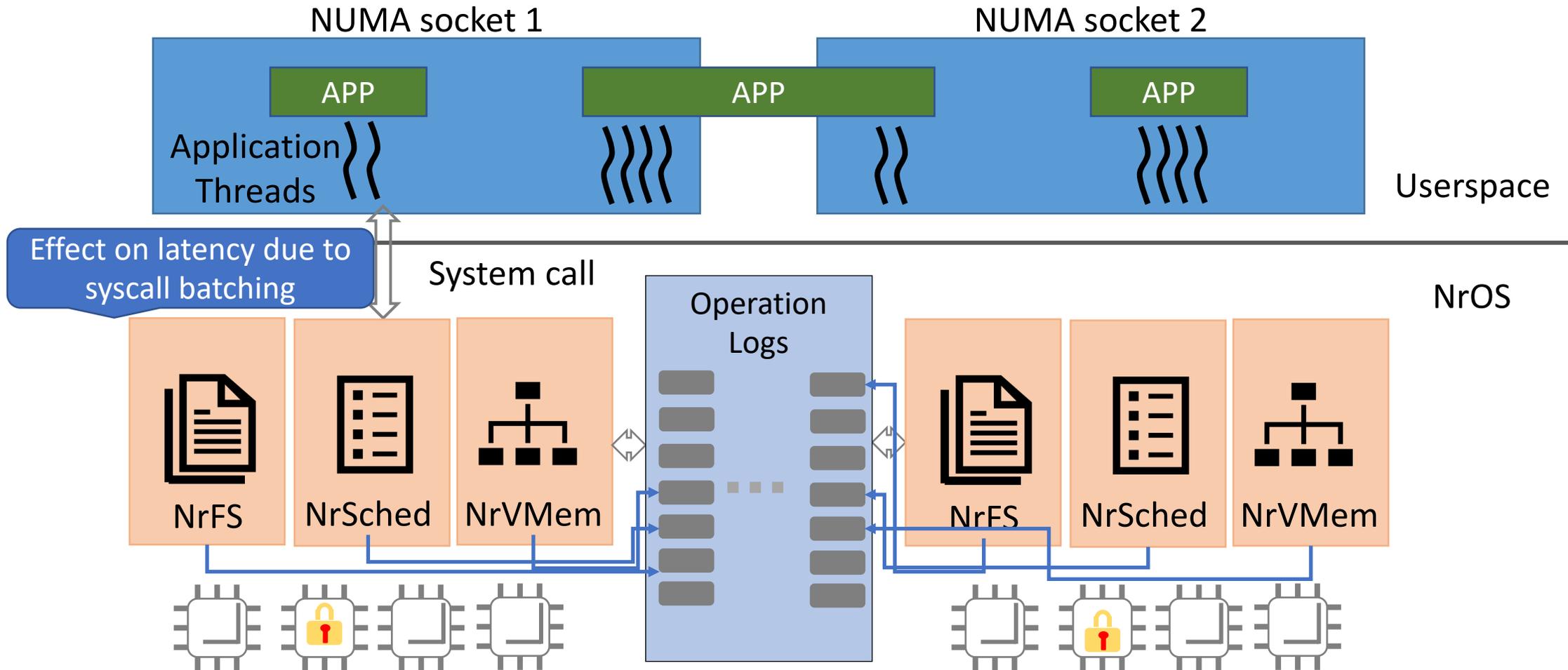
NrOS Summary



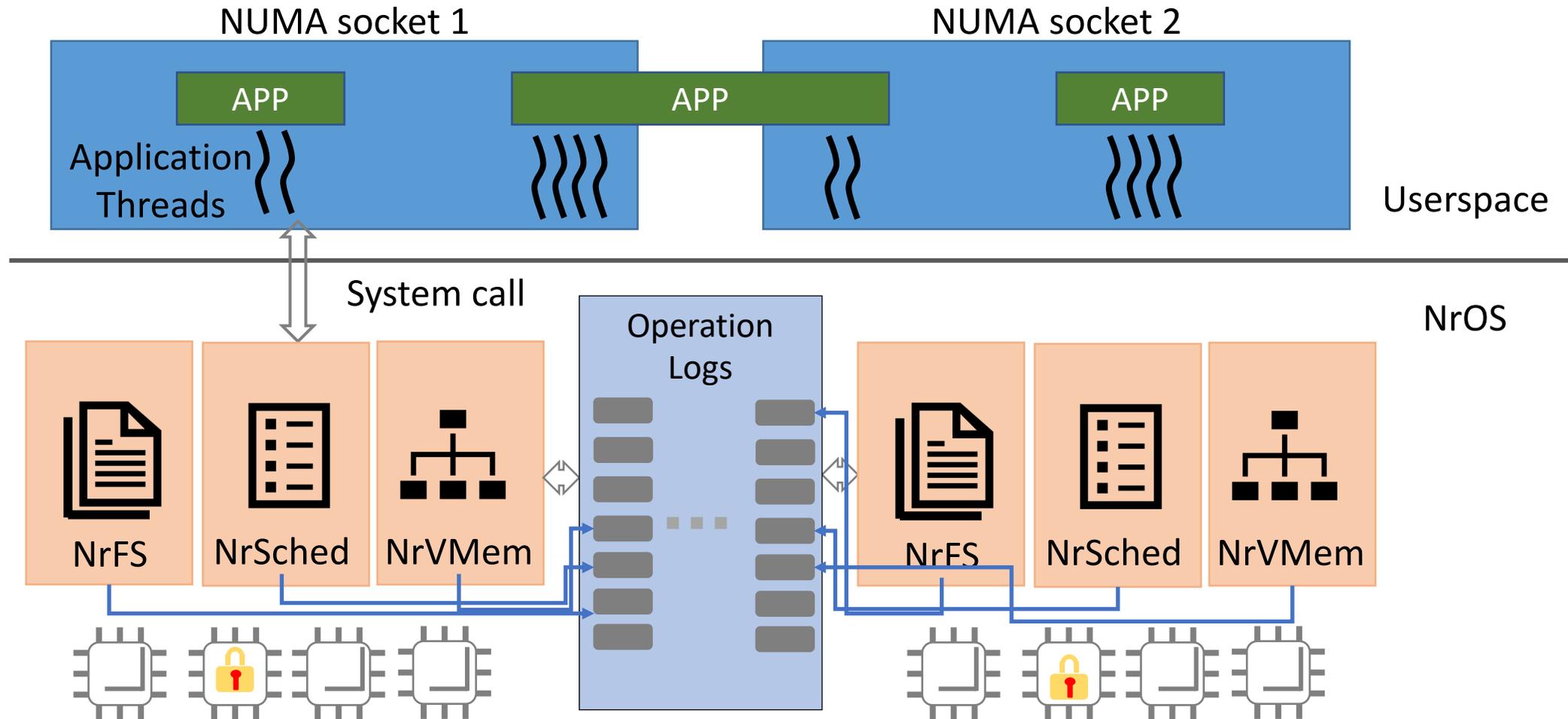
NrOS Summary



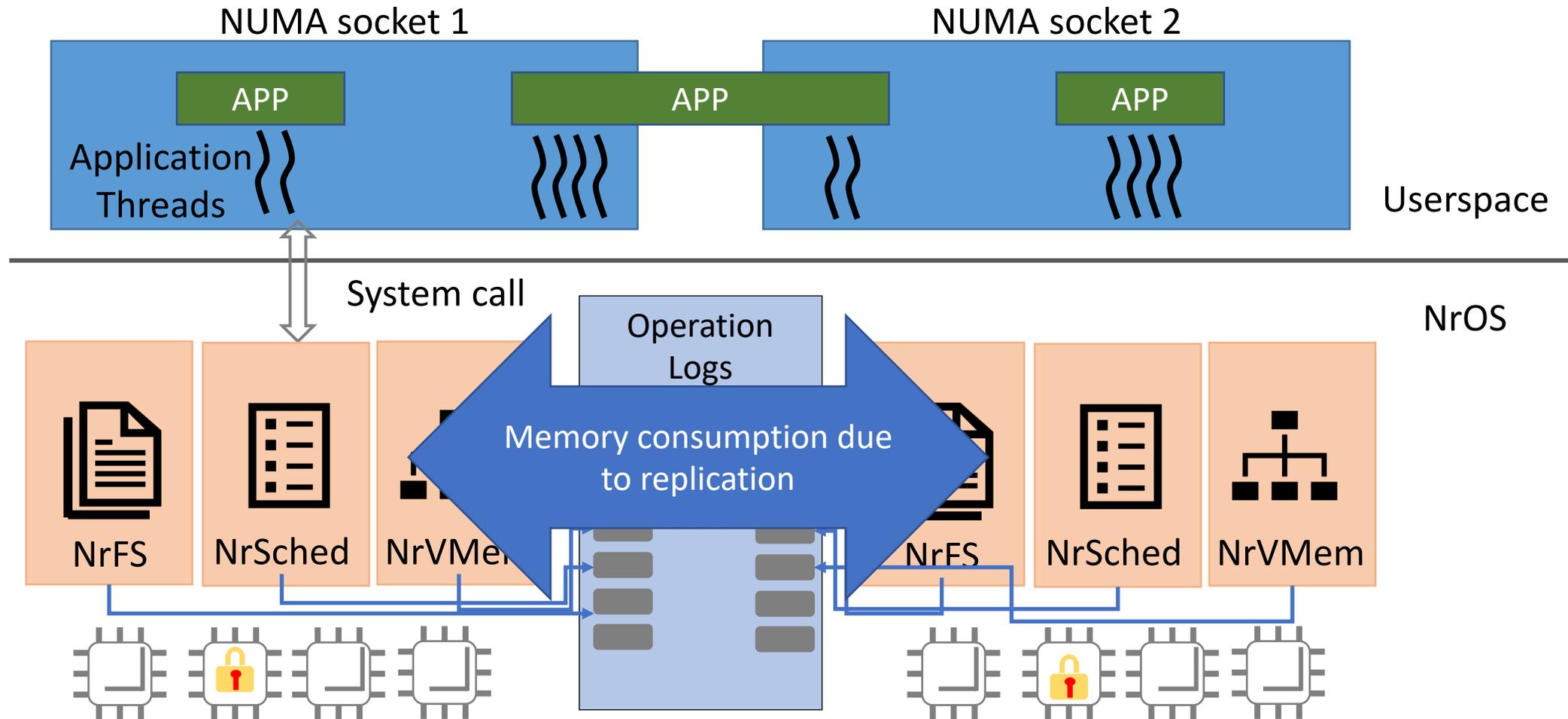
NrOS Summary



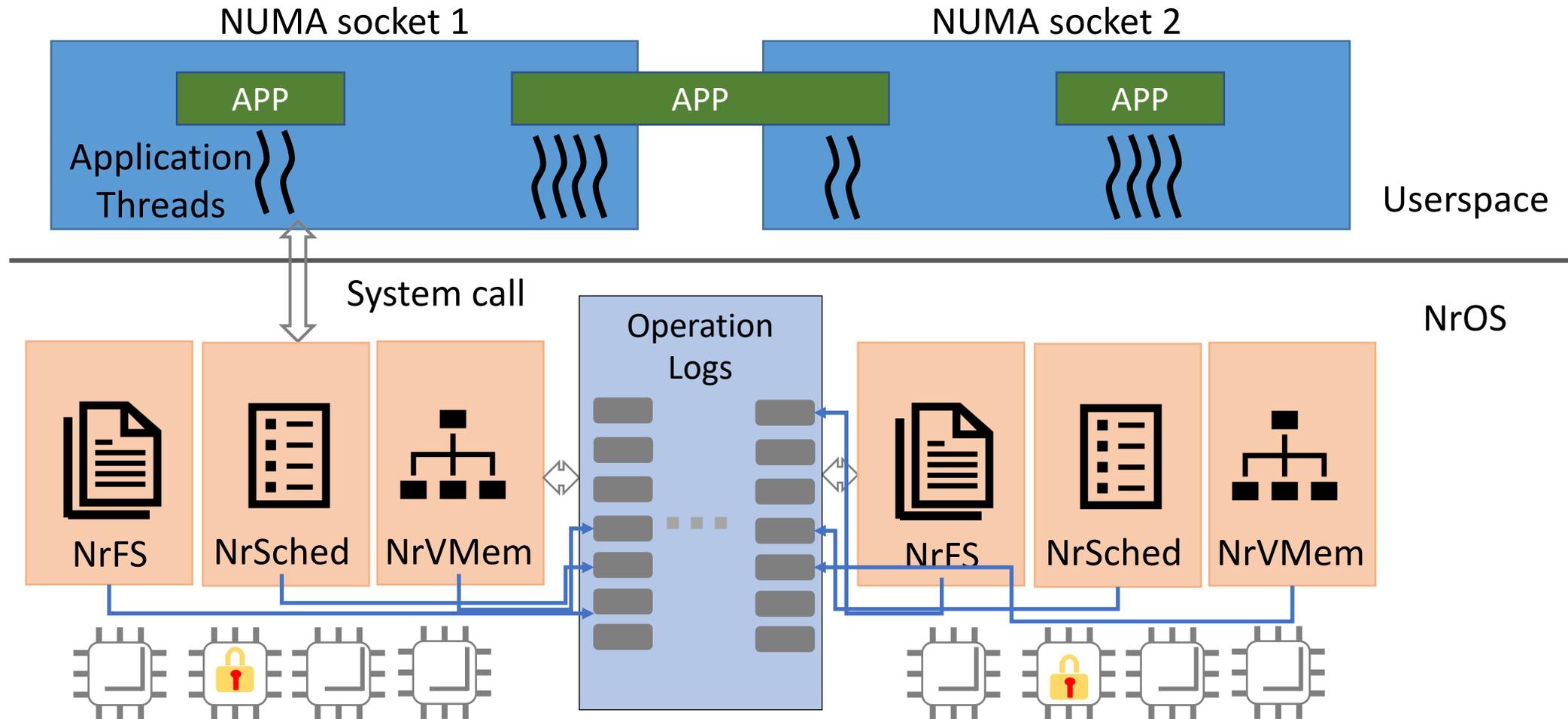
NrOS Summary



NrOS Summary



NrOS Summary



NRkernel: Key takeaways

- Simple, black-box approach (one size fits all) for concurrency
- Decomposes the problem of building scalable kernel data-structures for OS kernels
- Trade-off memory/CPU for performance: Replicate, make all reads NUMA local

<https://nrkernel.systems>