



**TECHNISCHE
UNIVERSITÄT
DRESDEN**

HPC - HIGH PERFORMANCE COMPUTING (SUPERCOMPUTING)

DISTRIBUTED OPERATING SYSTEMS, SCALABILITY, SS 2014

Hermann Härtig

Understand

- Systems Software for “High Performance Computing” (HPC), today & expected
- MPI as a common programming model
- What is “noise”?
- How to use incomplete information for informed decisions
- Advanced Load Balancing techniques (heuristics)

Characteristics of MPP Systems:

- Highly optimised interconnect networks
- Distributed memory
- Size today: few 100000 CPUs (cores) + XXL GPU

Successful Applications:

- CPU intensive computation, massively parallel Applications, small execution/communication ratios, weak and strong scaling
- Cloud ?

Not used for:

- Transaction-management systems
- Unix-Workstation + Servers

Characteristics of Cluster Systems:

- Use COTS (common off the shelf) PCs/Servers and COTS networks
- Size: No principle limits

Successful Applications:

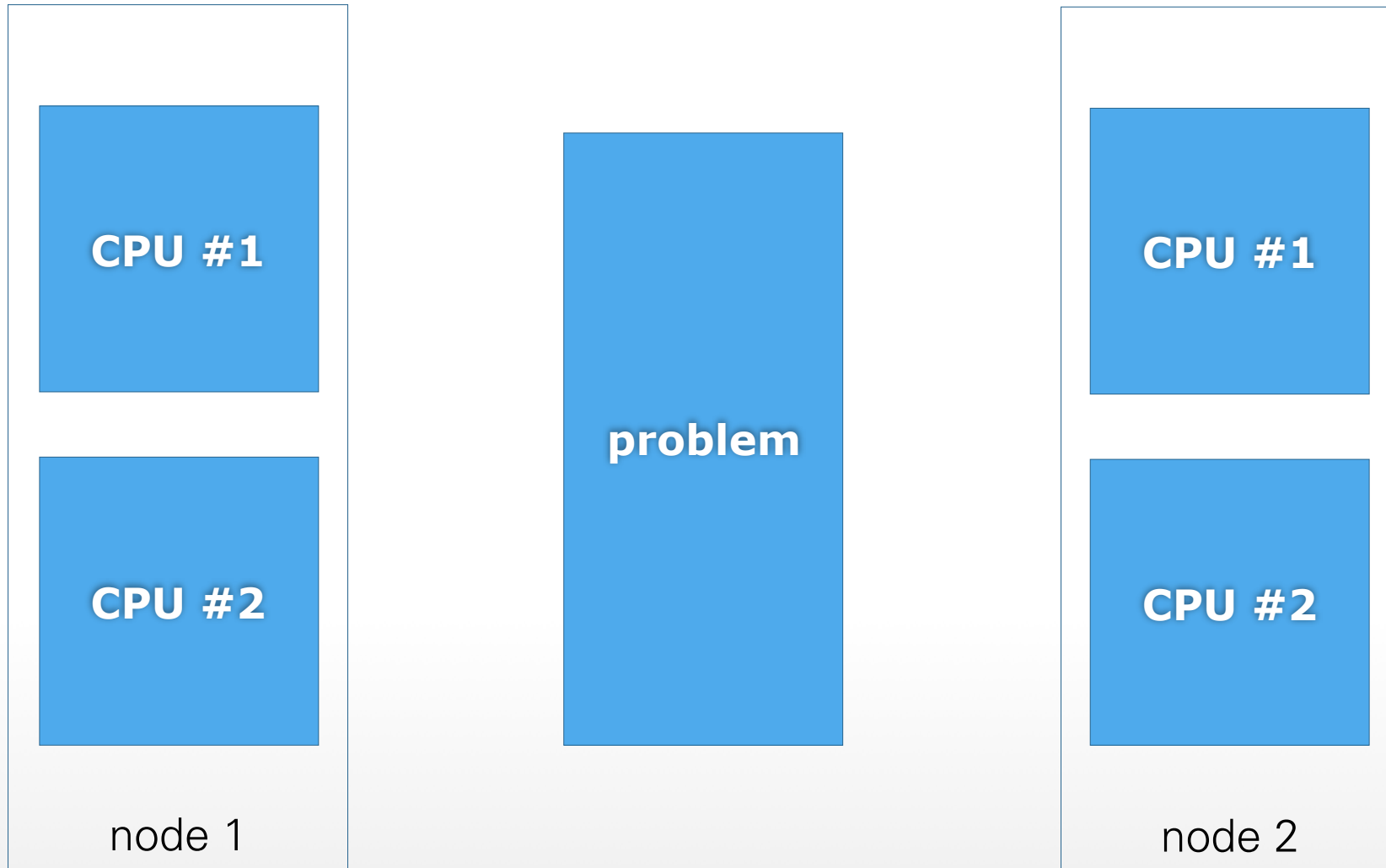
- CPU intensive computation, massively parallel Applications, larger execution/communication ratios, weak scaling
- Data Centers, google apps
- Cloud, Virtual Machines

Not used for:

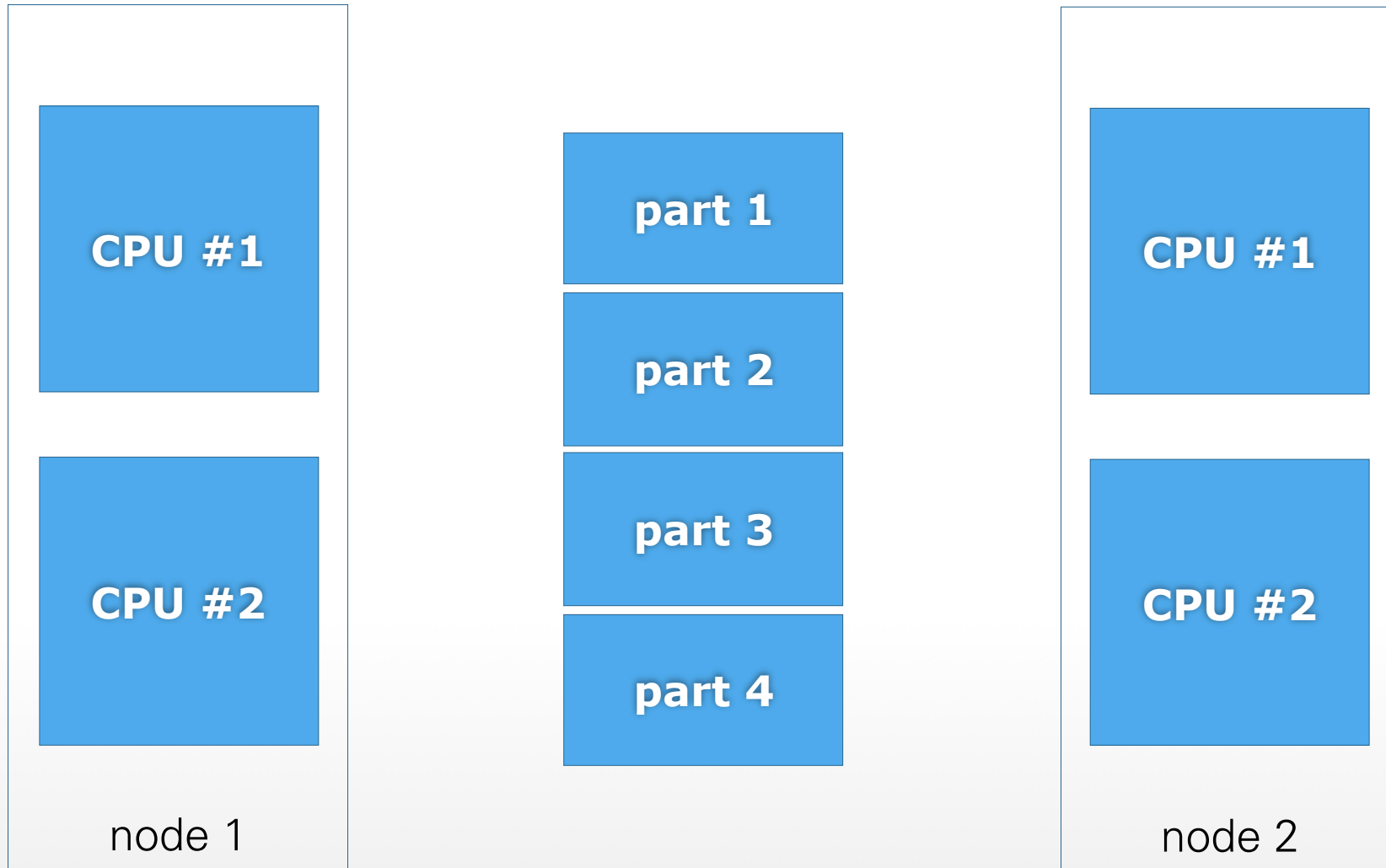
- Transaction-management system

- Michael Flynn (1966):
SISD, SIMD, MIMD, (MISD)SIMD
- SPMD: Single Program Multiple Data
Same program runs on “all” nodes
works on split-up data
asynchronously but with explicit synch points
implementations: message passing/shared
memory/...paradigms: “map/reduce” (google) /
GCD (apple) / task queues / ...
- often: while (true) { work; exchange data
(barrier)}

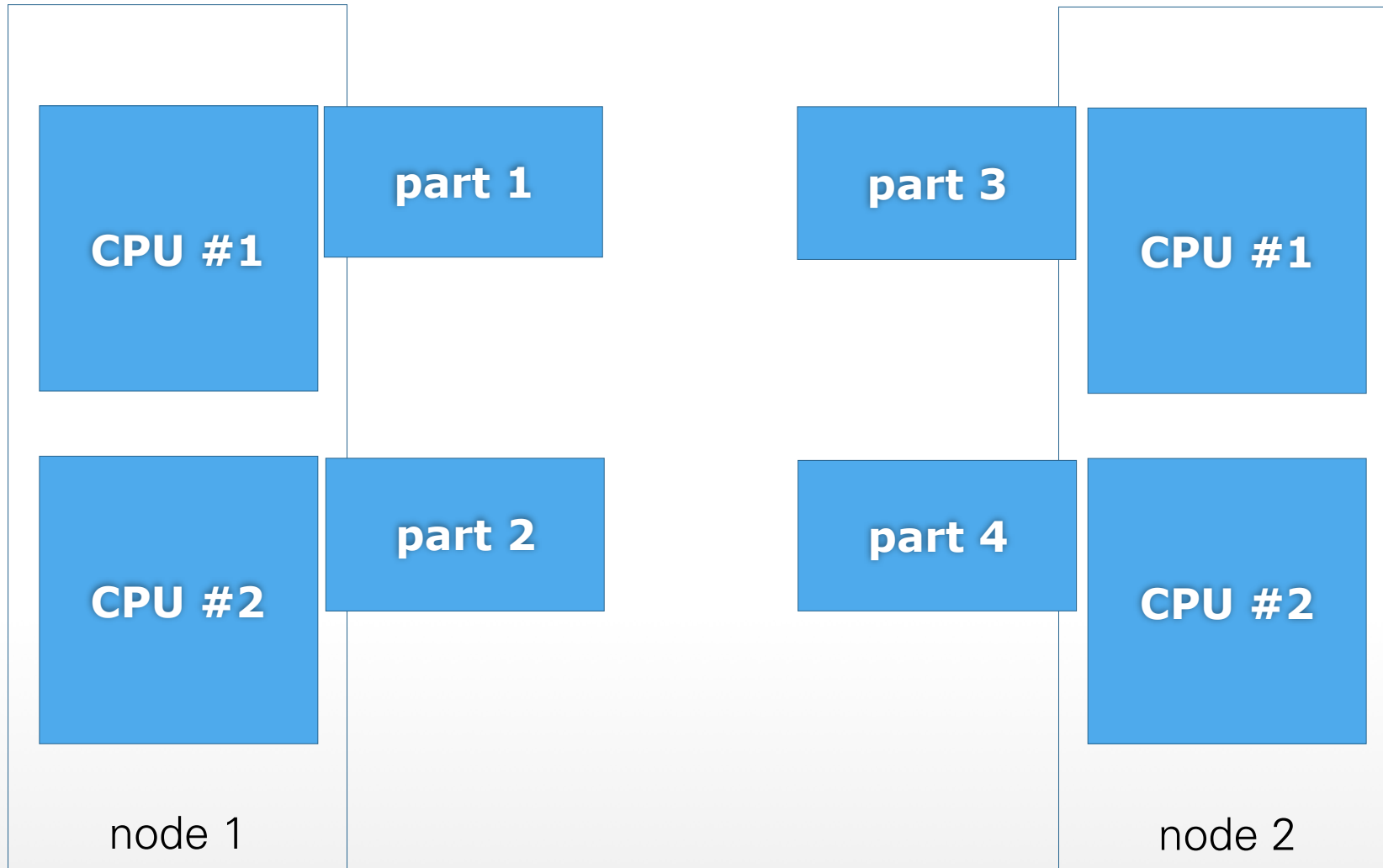
DIVIDE AND CONQUER



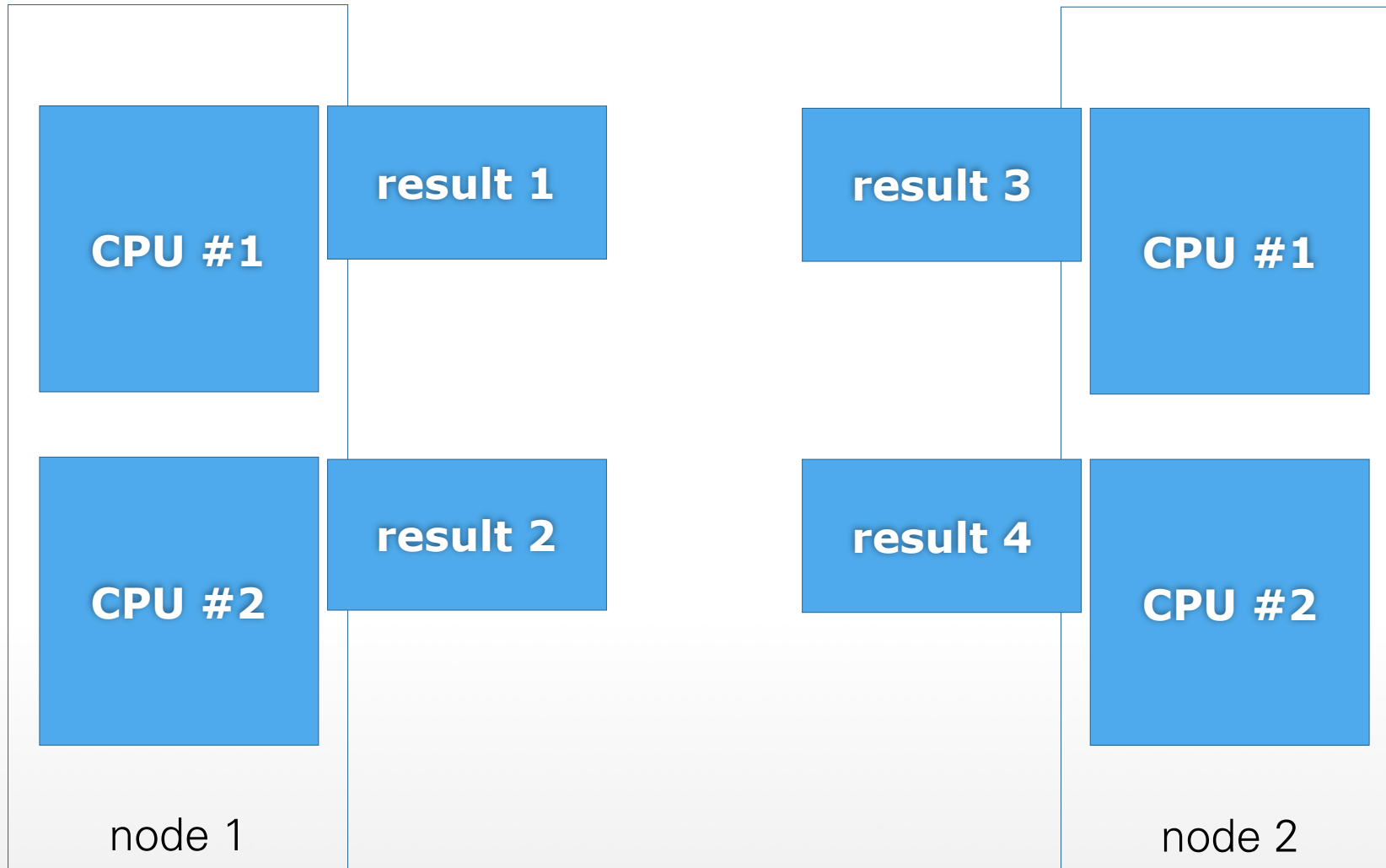
DIVIDE AND CONQUER



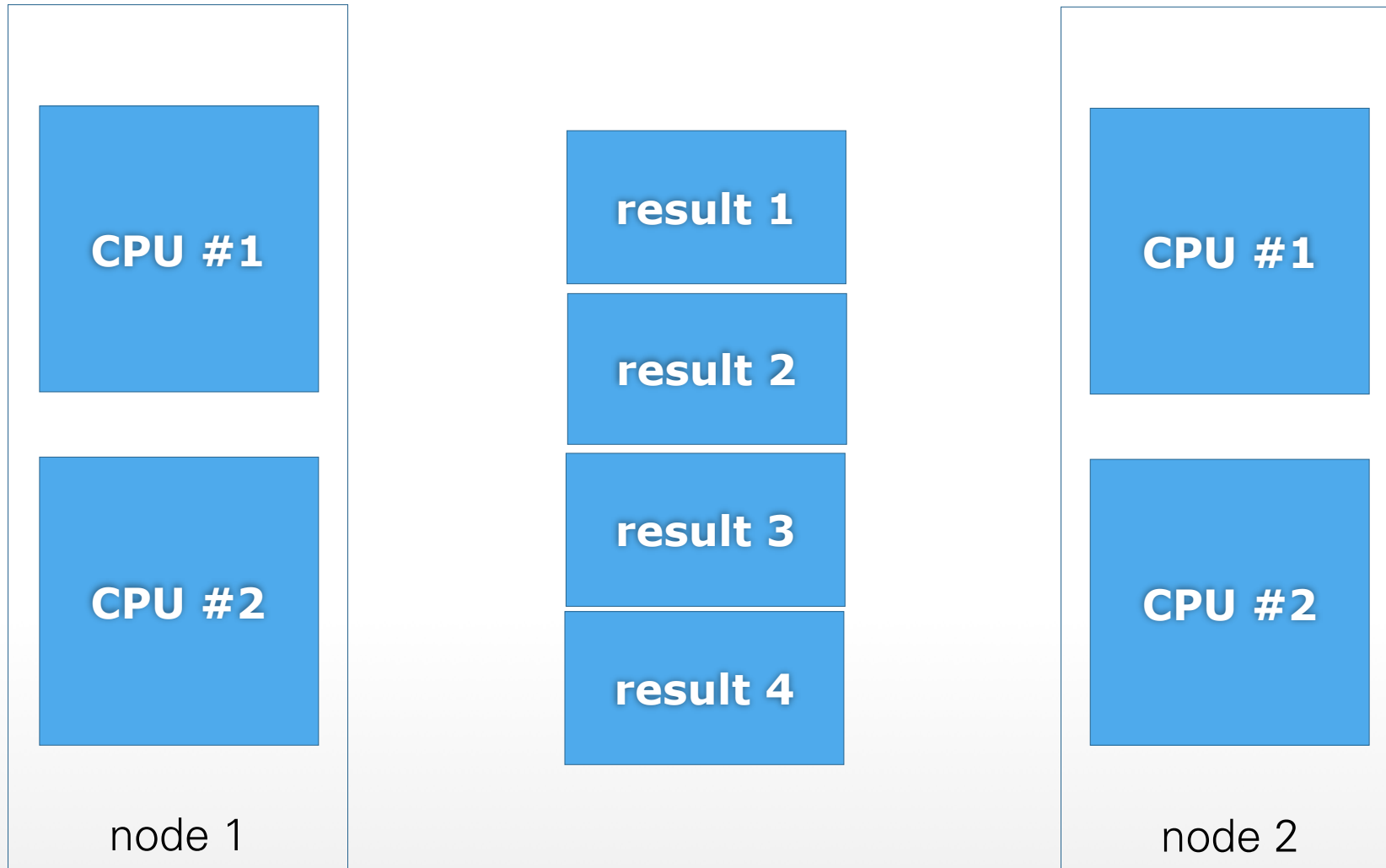
DIVIDE AND CONQUER



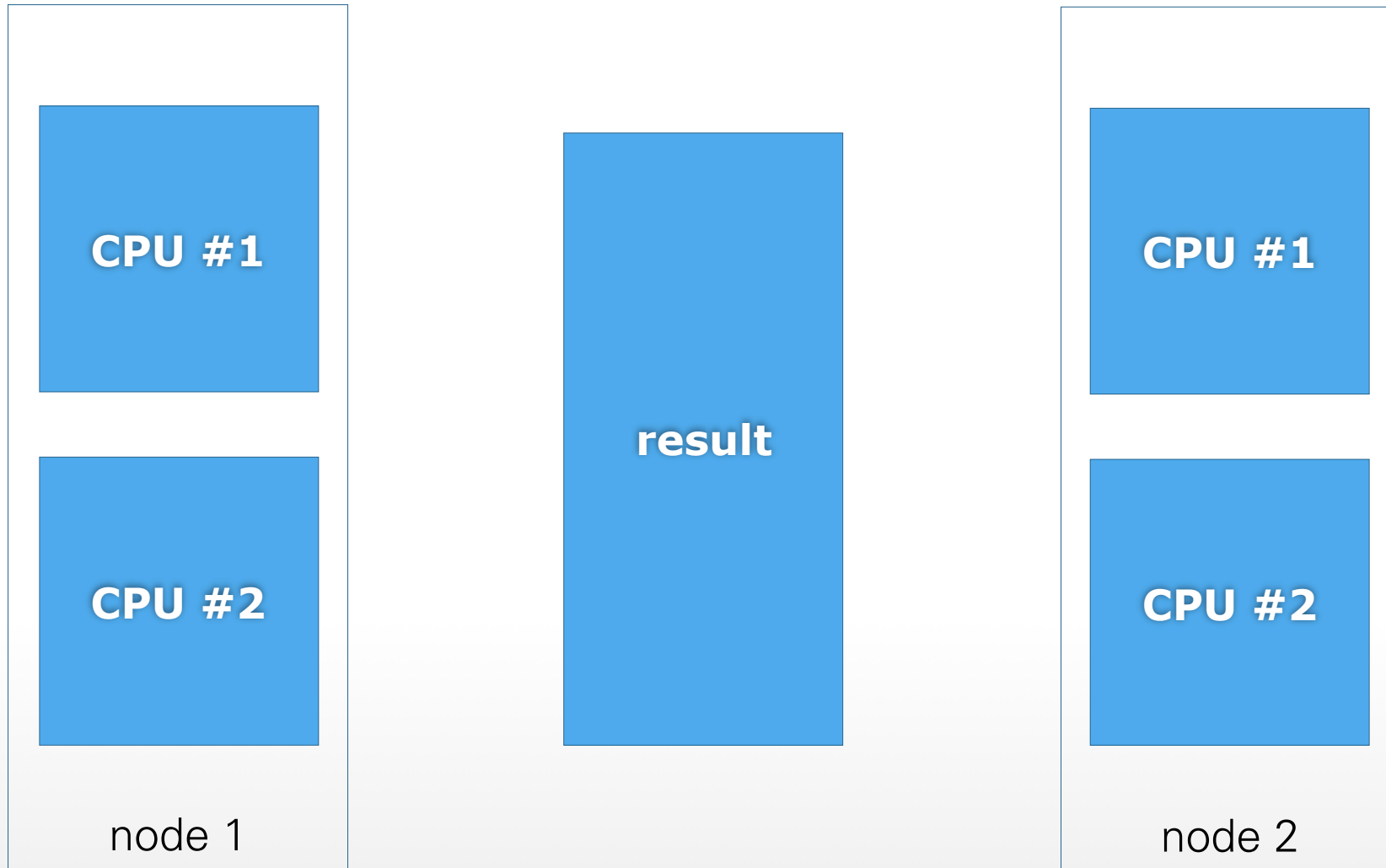
DIVIDE AND CONQUER

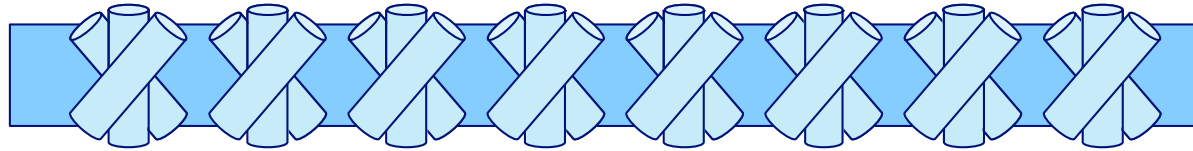


DIVIDE AND CONQUER

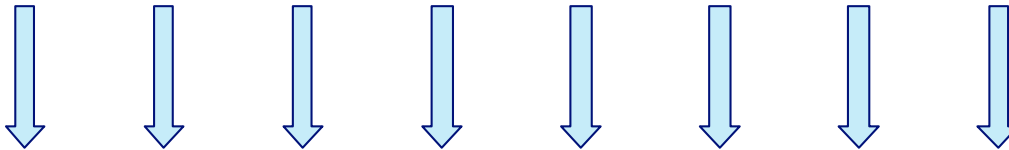


DIVIDE AND CONQUER

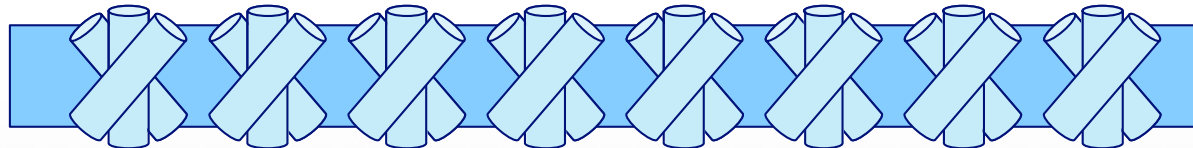




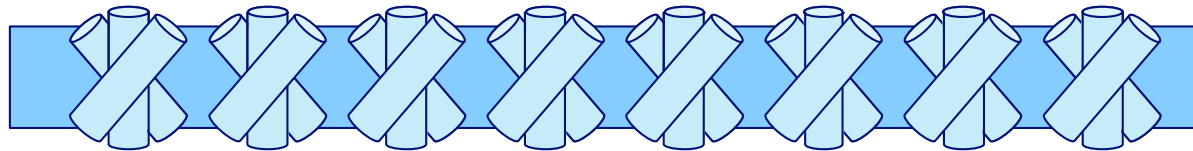
Communication



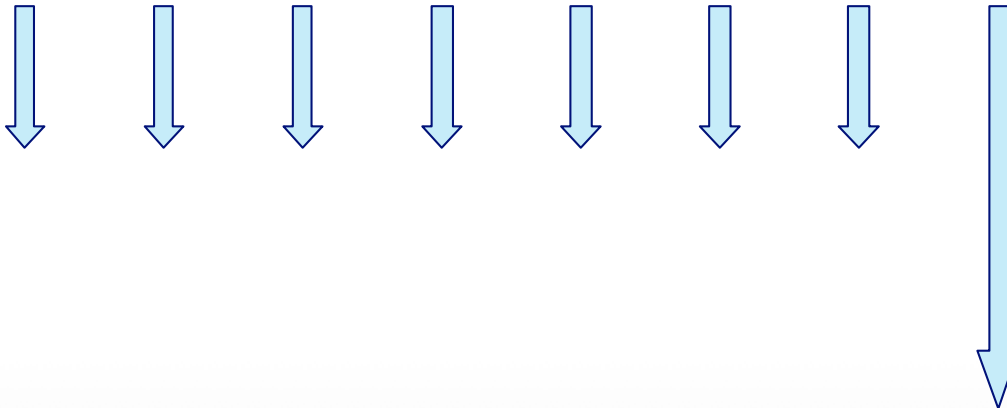
Computation



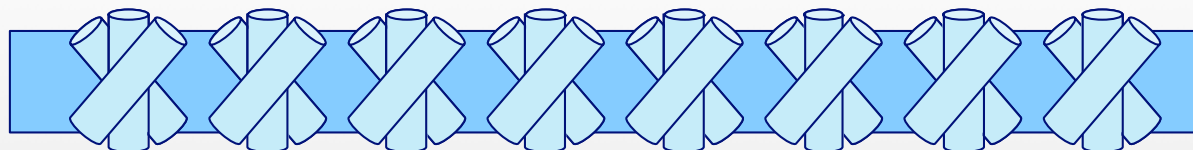
Communication



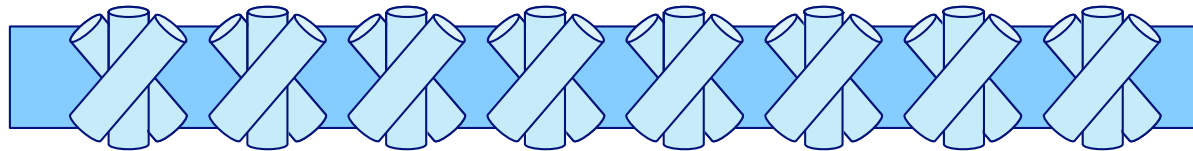
Communication



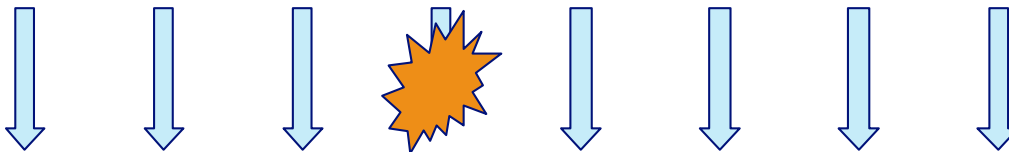
Computation



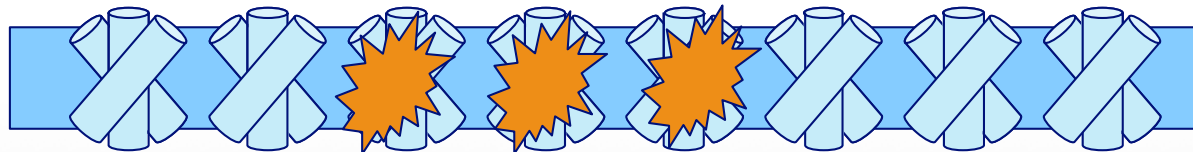
Communication



Communication



Computation



Communication

Compute; communicate; compute; ...

- Examples (idealized, take with grain of salt !!!):
 - Compute: 10 micro, 100 micro, 1 ms
 - Communicate: 5 micro, 10 micro, 100 micro, 1ms
assuming here: communication cannot be sped up

Amdahl's law: $1 / (1 - P + P/N)$

- P: section that can be parallelized
- 1-P: serial section
- N: number of CPUs

Compute(= parallel section),
communicate(= serial section)



possible speedup for $N = \infty$

- 1ms, 100 μ s: 1/0.1 → 10
- 1ms, 1 μ s: 1/0.001 → 1000
- 10 μ s, 1 μ s: 0.01/0.001 → 10
- ...

Strong:

- accelerate same problem size

Weak:

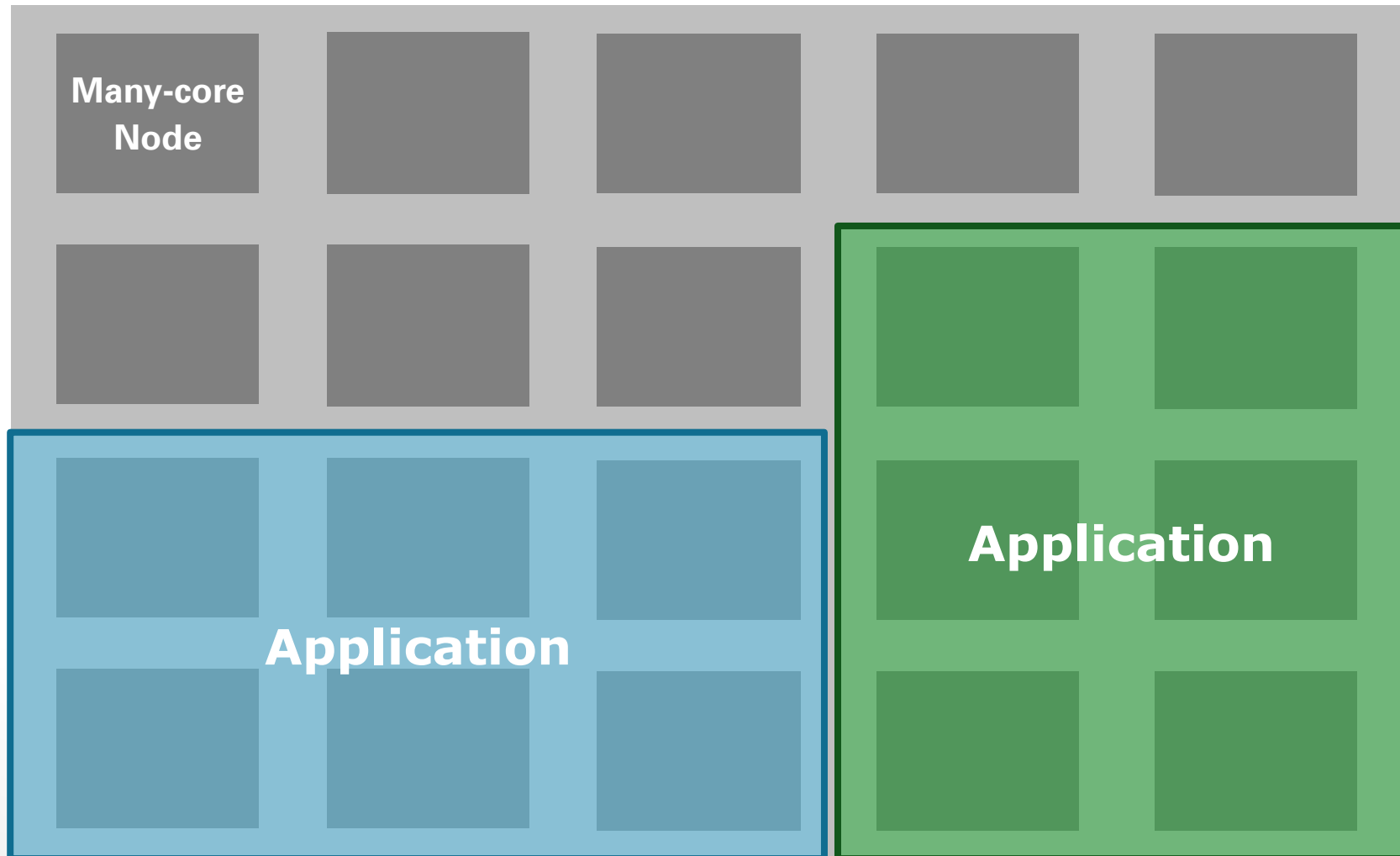
- extend to larger problem size

Jitter, "Noise", "micro scrabblers":

- Occasional addition to computation/communication time in one or more processes
- Holds up all other processes

Compute(= parallel section),
jitter (→ add to serial section),
communicate(= serial section):
possible speedup for $N=\infty$

- 1ms, 100 μ s, 100 μ s: 1/0.2 → 5 (10)
- 1ms, 100 μ s, 1 μ s: 1/0.101 → 10 (1000)
- 10 μ s, 10 μ s, 1 μ s: 0.01/0.011 → 1 (10)

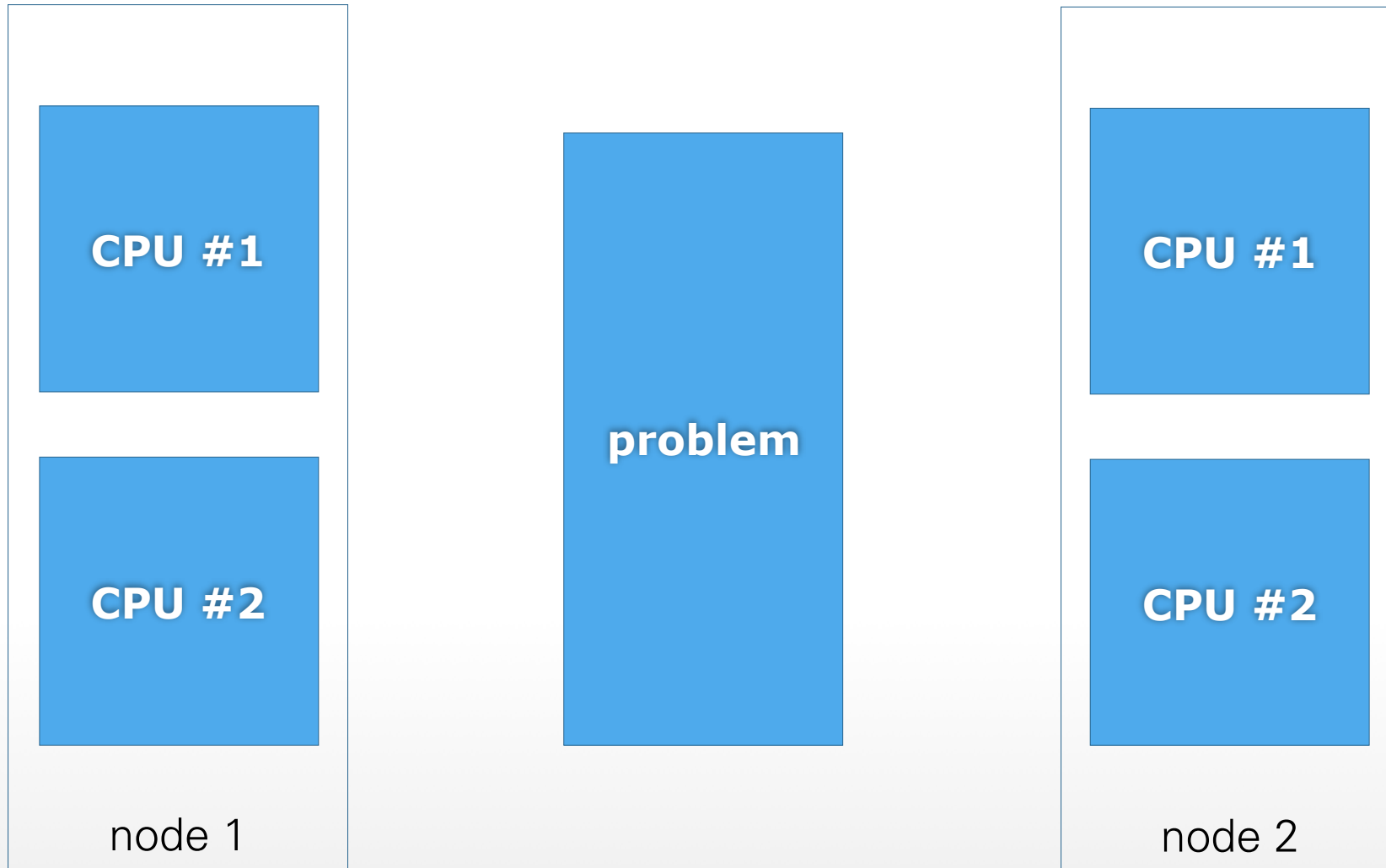


- dedicate full partition to application
(variant: “gang scheduling”)
- load balancing done (tried) by applications or
user-level runtime (Charm++)
- avoid OS calls
- “scheduler”:
manages queue of application processes
assigns partitions to applications
supervises run-time
- applications run from checkpoint to checkpoint

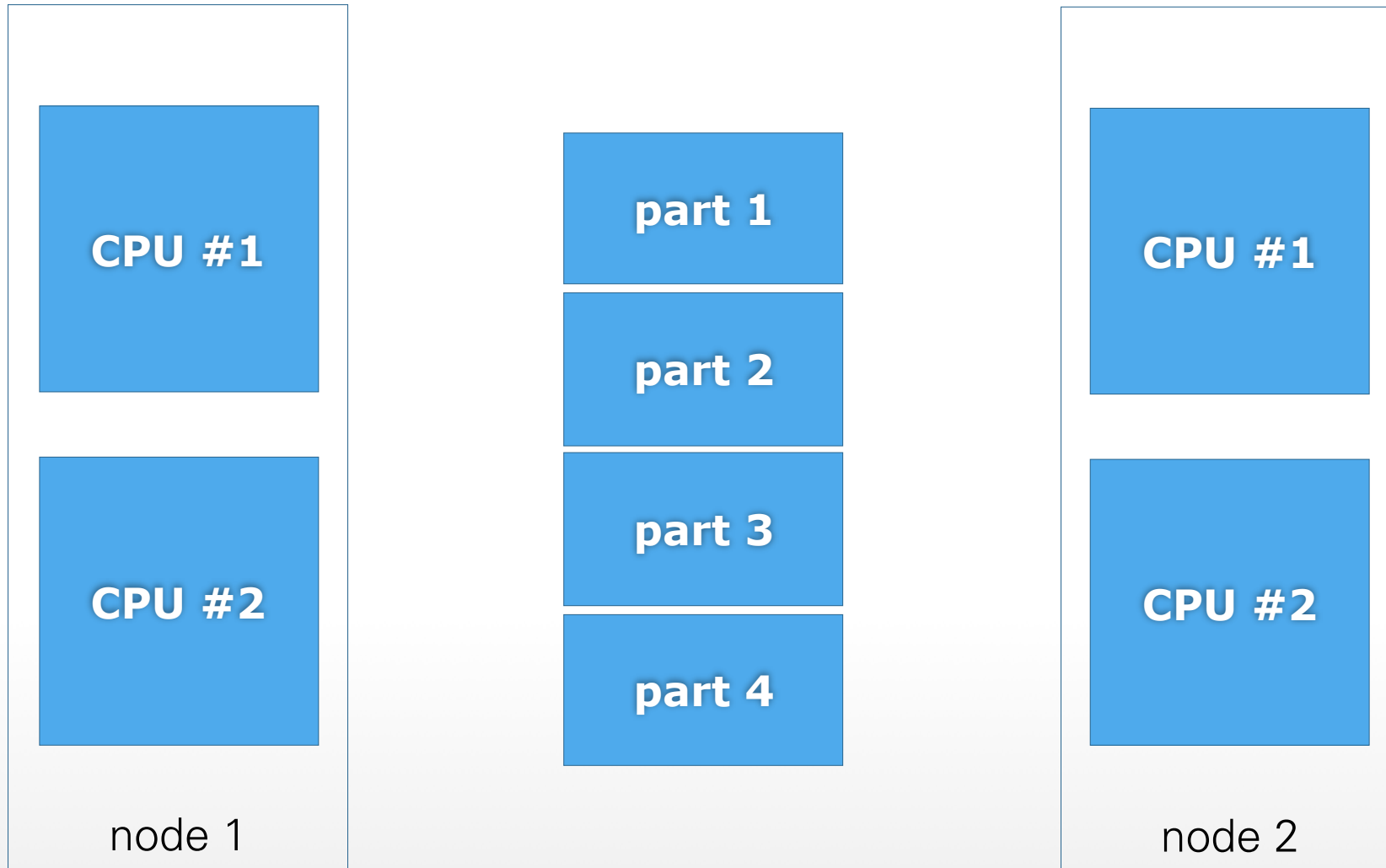
- nodes access remote memory via load/store operations
- busy waiting across nodes (within partition)
- barrier ops supported by network
- compare&exchange on remote memory operation
- no OS calls for message ops (busy waiting)

- Library for message-oriented parallel programming
- Programming model:
 - Multiple instances of same program
 - Independent calculation
 - Communication, synchronization

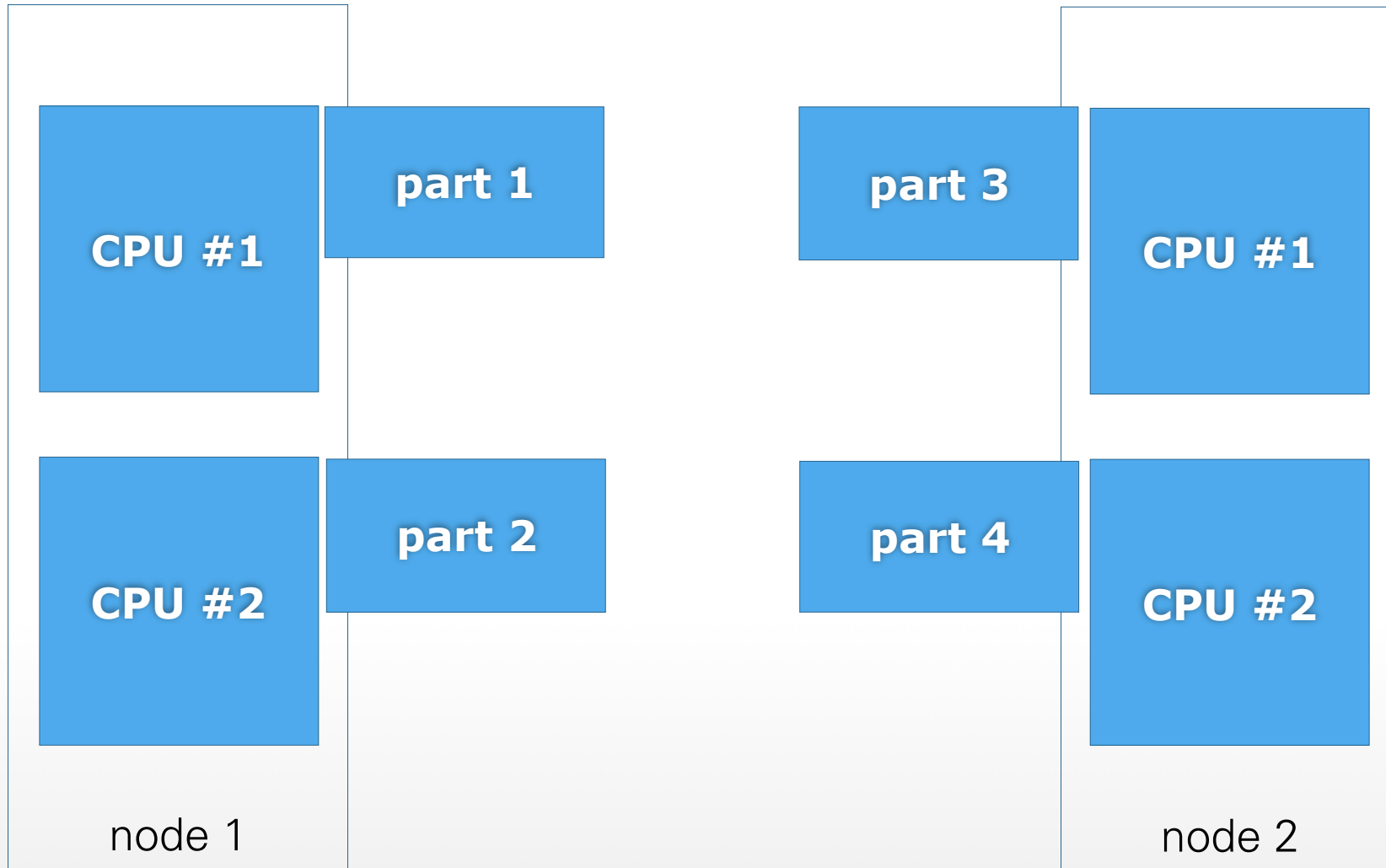
DIVIDE AND CONQUER



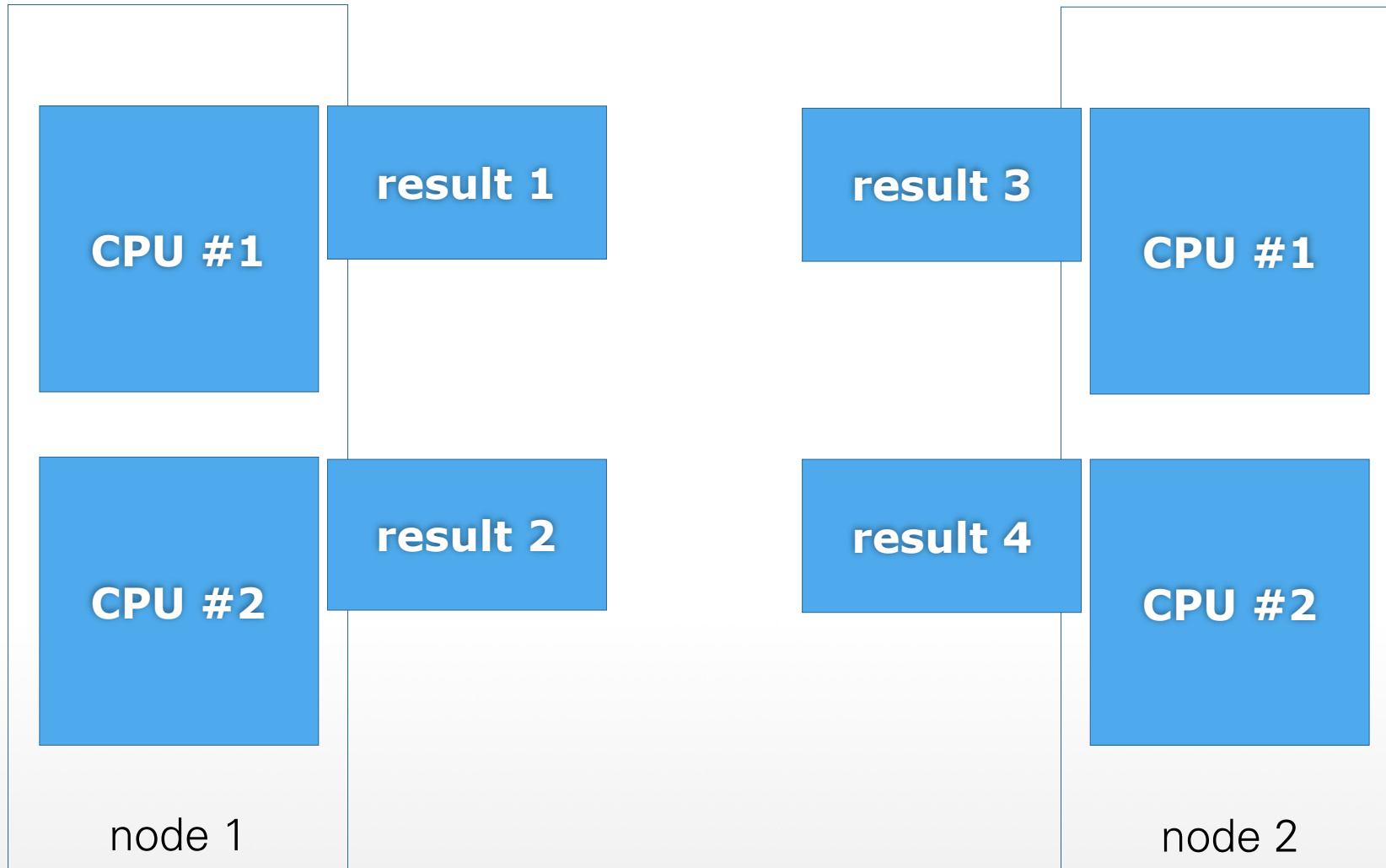
DIVIDE AND CONQUER



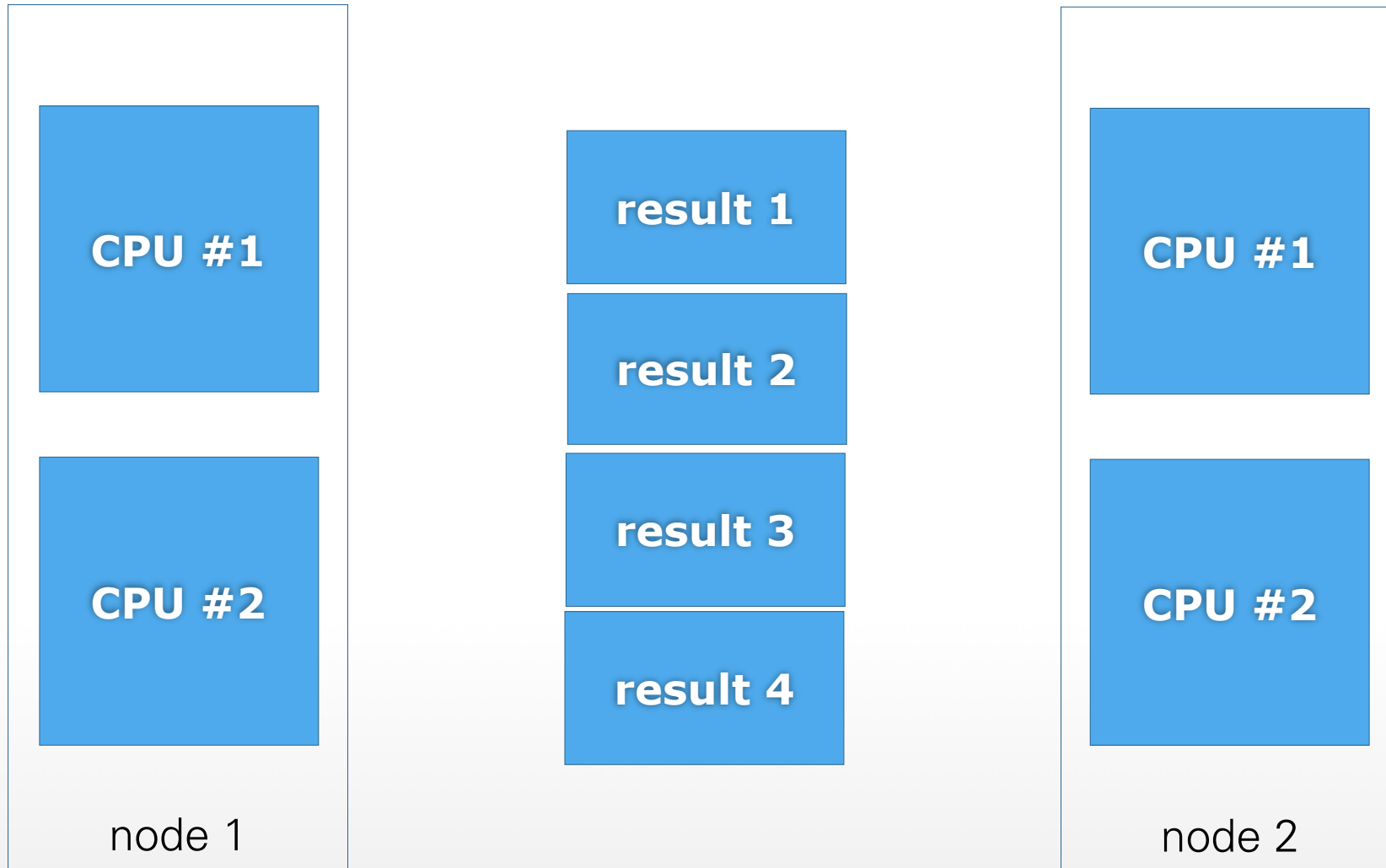
DIVIDE AND CONQUER



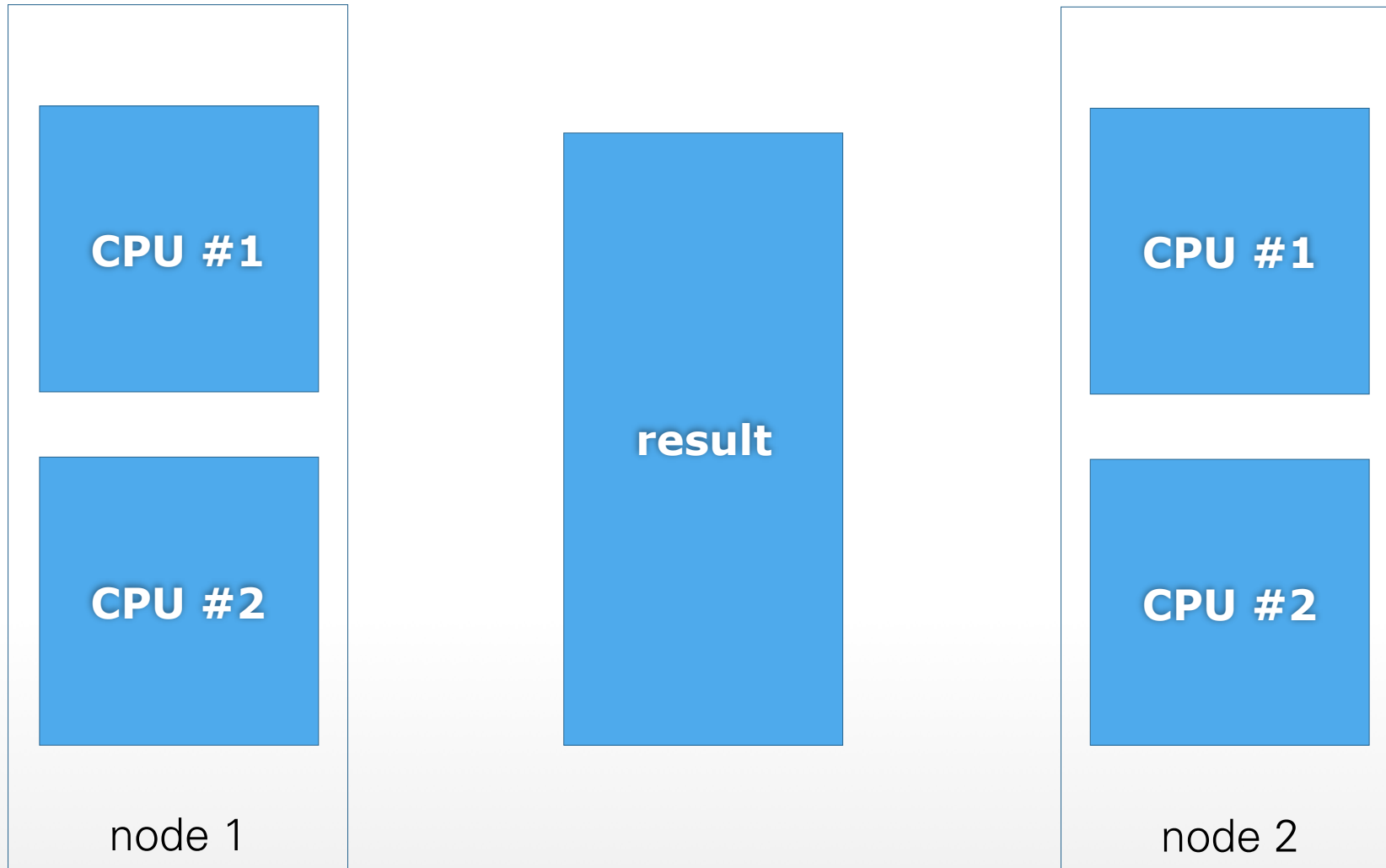
DIVIDE AND CONQUER



DIVIDE AND CONQUER



DIVIDE AND CONQUER



- MPI program is started on all processors
- `MPI_Init()`, `MPI_Finalize()`
- Communicators (e.g., `MPI_COMM_WORLD`)
 - `MPI_Comm_size()`
 - `MPI_Comm_rank()`: “Rank” of process within this set
 - Typed messages
- Dynamically create and spread processes using `MPI_Spawn()` (since MPI-2)

- Communication

- Communication
 - Point-to-point

```
MPI_Send(  
    void* buf,  
    int count,  
    MPI_Datatype,  
    int dest,  
    int tag,  
    MPI_Comm comm  
)
```


- Communication
 - Point-to-point

```
MPI_Recv(  
    void* buf,  
    int count,  
    MPI_Datatype,  
    int source,  
    int tag,  
    MPI_Comm comm,  
    MPI_Status *status  
)
```

- Communication
 - Point-to-point
 - Collectives

```
MPI_Bcast(  
    void* buffer,  
    int count,  
    MPI_Datatype,  
    int root,  
    MPI_Comm comm  
)
```

- Communication
 - Point-to-point
 - Collectives

```
MPI_Reduce(  
    void* sendbuf,  
    void *recvbuf,  
    int count  
    MPI_Datatype,  
    MPI_Op op,  
    int root,  
    MPI_Comm comm  
)
```

- Communication
 - Point-to-point
 - Collectives
- Synchronization

- Communication
 - Point-to-point
 - Collectives
- Synchronization
 - Test

```
MPI_Test(  
    MPI_Request* request,  
    int *flag,  
    MPI_Status *status  
)
```

- Communication
 - Point-to-point
 - Collectives
- Synchronization
 - Test
 - Wait

```
MPI_Wait(  
    MPI_Request* request,  
    MPI_Status *status  
)
```

- Communication
 - Point-to-point
 - Collectives
- Synchronization
 - Test
 - Wait
 - Barrier

```
MPI_Barrier(  
    MPI_Comm comm  
)
```

BLOCK AND SYNC

	blocking call	non-blocking call
synchronous communication		
asynchronous communication		

	blocking call	non-blocking call
synchronous communication	returns when message has been delivered	
asynchronous communication		

	blocking call	non-blocking call
synchronous communication	returns when message has been delivered	
asynchronous communication	returns when send buffer can be reused	

	blocking call	non-blocking call
synchronous communication	returns when message has been delivered	returns immediately, following test/wait checks for delivery
asynchronous communication	returns when send buffer can be reused	

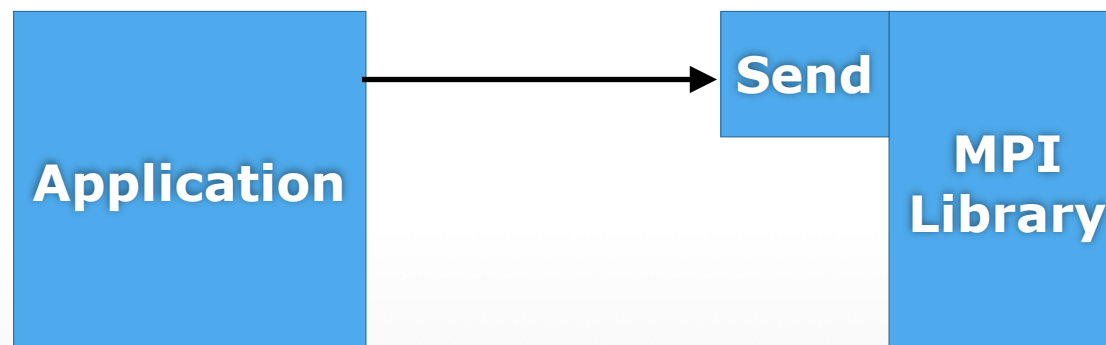
	blocking call	non-blocking call
synchronous communication	returns when message has been delivered	returns immediately, following test/wait checks for delivery
asynchronous communication	returns when send buffer can be reused	returns immediately, following test/wait checks for send buffer

```
int rank, total;
MPI_Init();
MPI_Comm_rank(MPI_COMM_WORLD, &rank);
MPI_Comm_size(MPI_COMM_WORLD, &total);

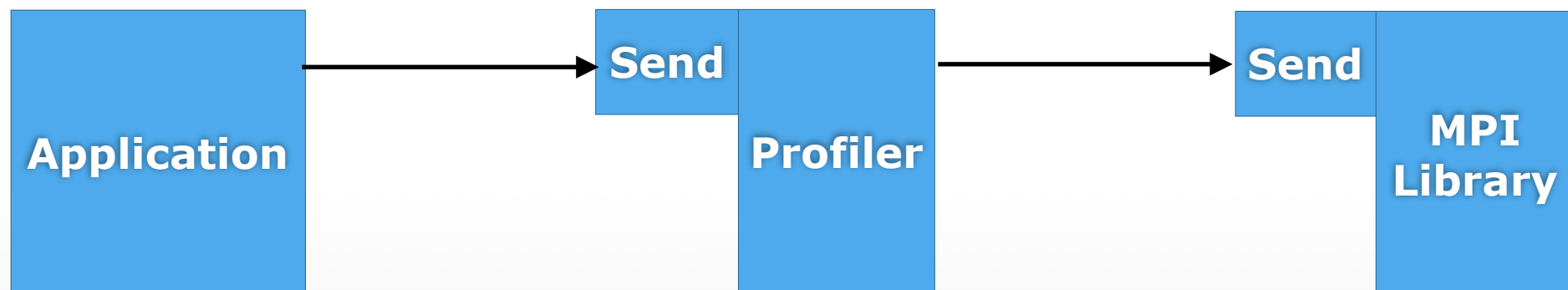
MPI_Bcast(...);
/* work on own part, determined by rank */

if (id == 0) {
    for (int rr = 1; rr < total; ++rr)
        MPI_Recv(...);
    /* Generate final result */
} else {
    MPI_Send(...);
}
MPI_Finalize();
```

- Interposition layer between library and application
- Originally designed for profiling



- Interposition layer between library and application
- Originally designed for profiling



- Large number of nodes:
 - Many compute cores
 - 1 or 2 service cores
- Failure rate exceeds checkpoint rate
- Fast local persistent storage on each node
- Not all cores available all the time (dark silicon due to heat/energy issues)
- Compute + communication heavy applications, may not be balanced
- short term changes of frequency ?

- for applications with extreme (bad) computation/
communications ratio:
NOT MUCH, but
-> avoid "noise", use common sense
- all others:
handle faults
use dark silicon
balance load
gossip
over decomposition & over subscription
predict execution times
use scheduling tricks
optimise for network/memory topology

Use common sense to avoid:

- OS usually not directly on the critical path, BUT OS controls: interference via interrupts, caches, network, memory bus, (RTS techniques)
- avoid or encapsulate side activities
- small critical sections (if any)
- partition networks to isolate traffic of different applications (HW: Blue Gene)
- do not run Python scripts or printer daemons in parallel

+ Hebrew Uni (Mosix team) + ZIB (FS team)
Fast and Fault-Tolerant Microkernel-based OS

- get rid of partitions
- use a micro-kernel (L4)
- OS supported load balancing
- use RAID for fast checkpoints

DFG-supported

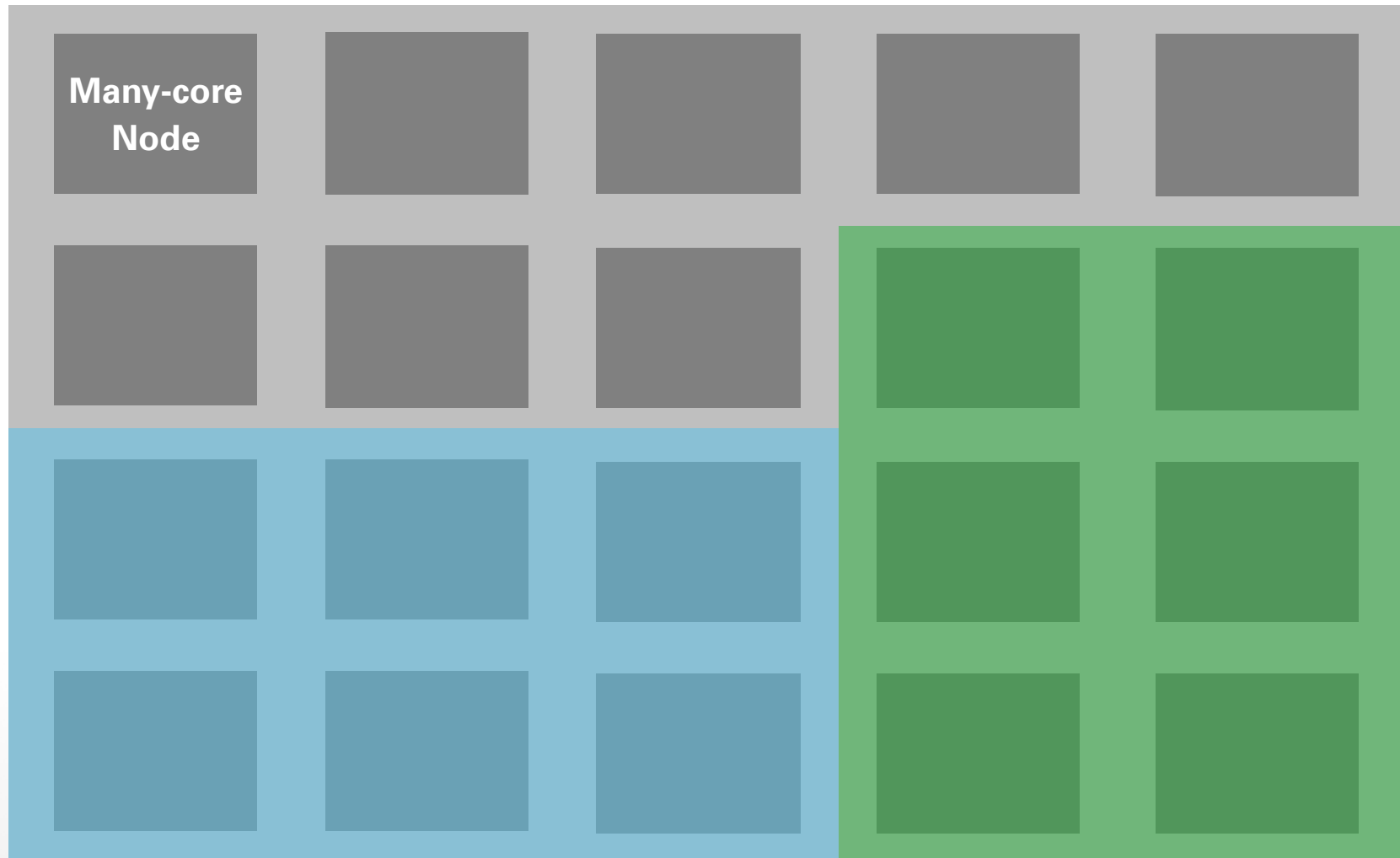
Microkernels, virtualization, split architectures

MOSIX-style online system management (gossip)

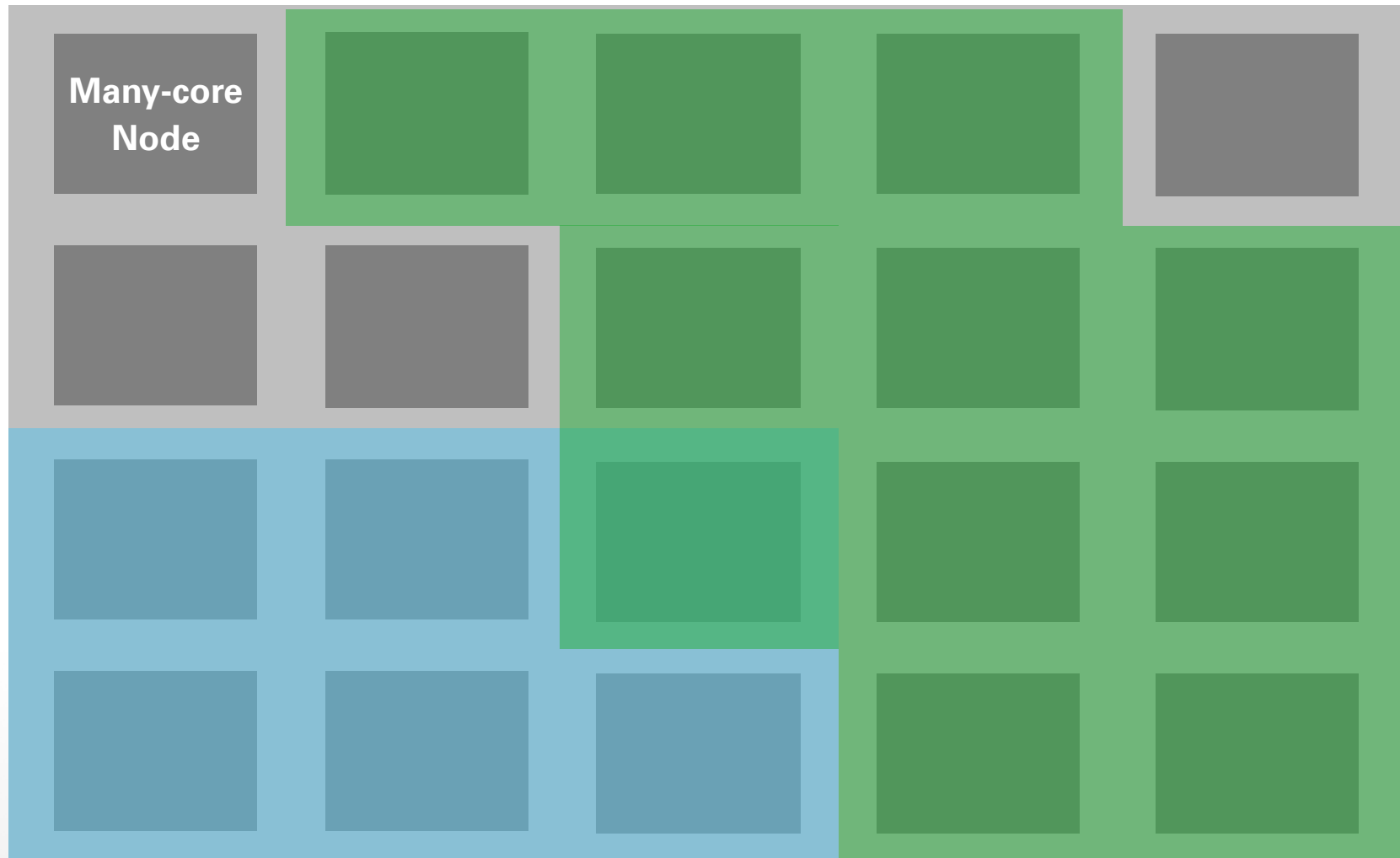
Distributed in-memory (on-node) checkpointing

MPI + applications

GOAL FOR EXASCALE HPC



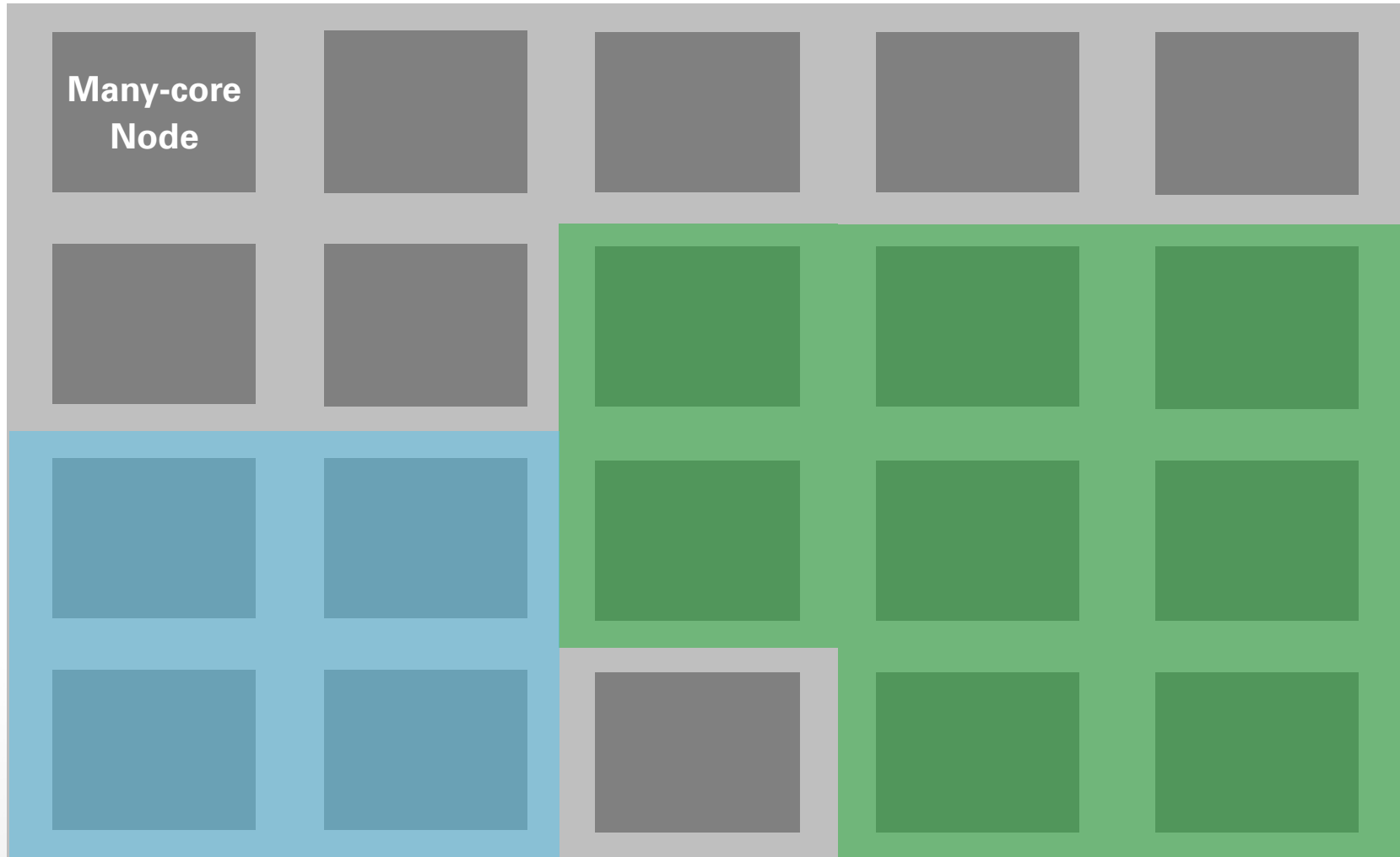
GOAL FOR EXASCALE HPC

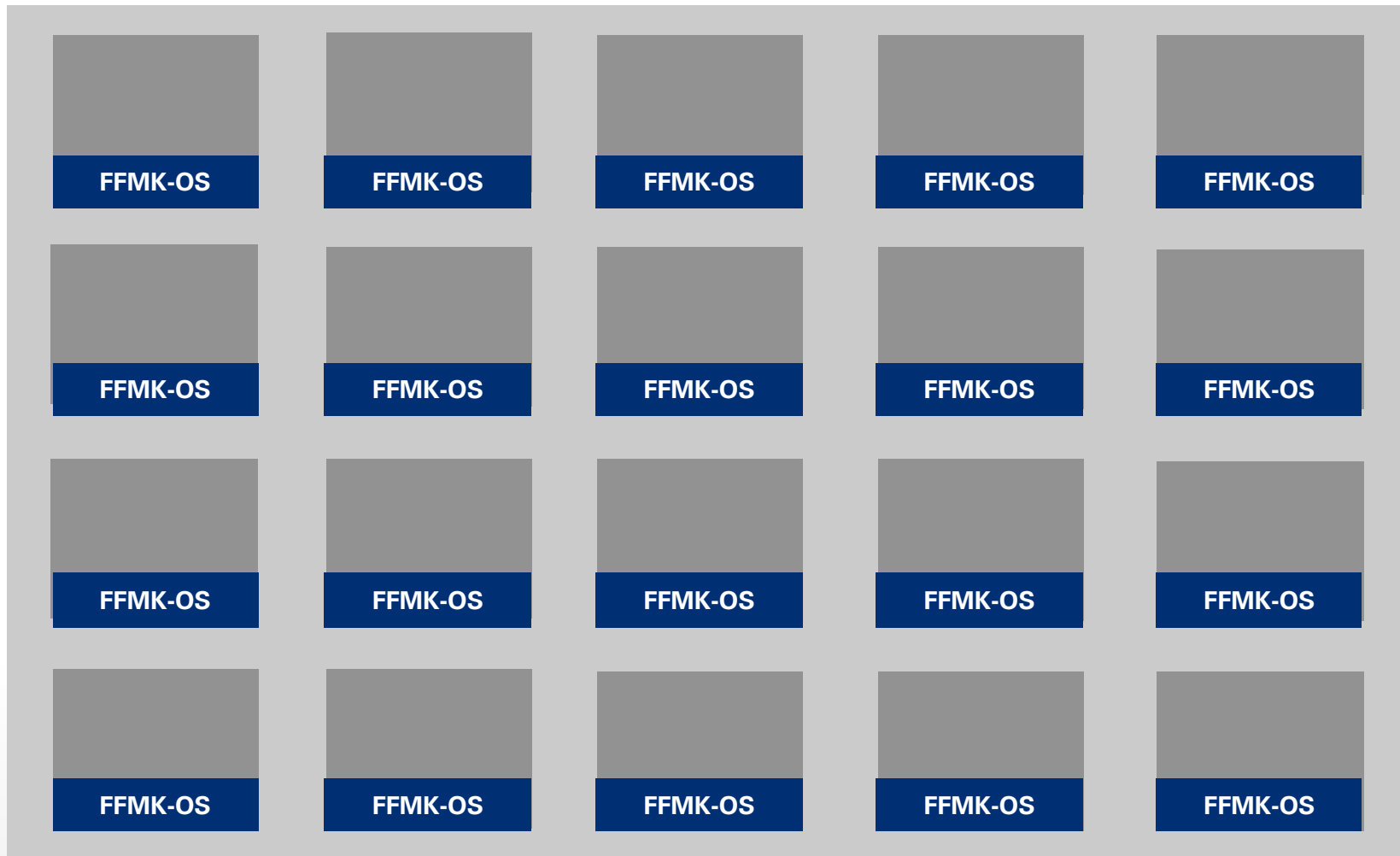


GOAL FOR EXASCALE HPC

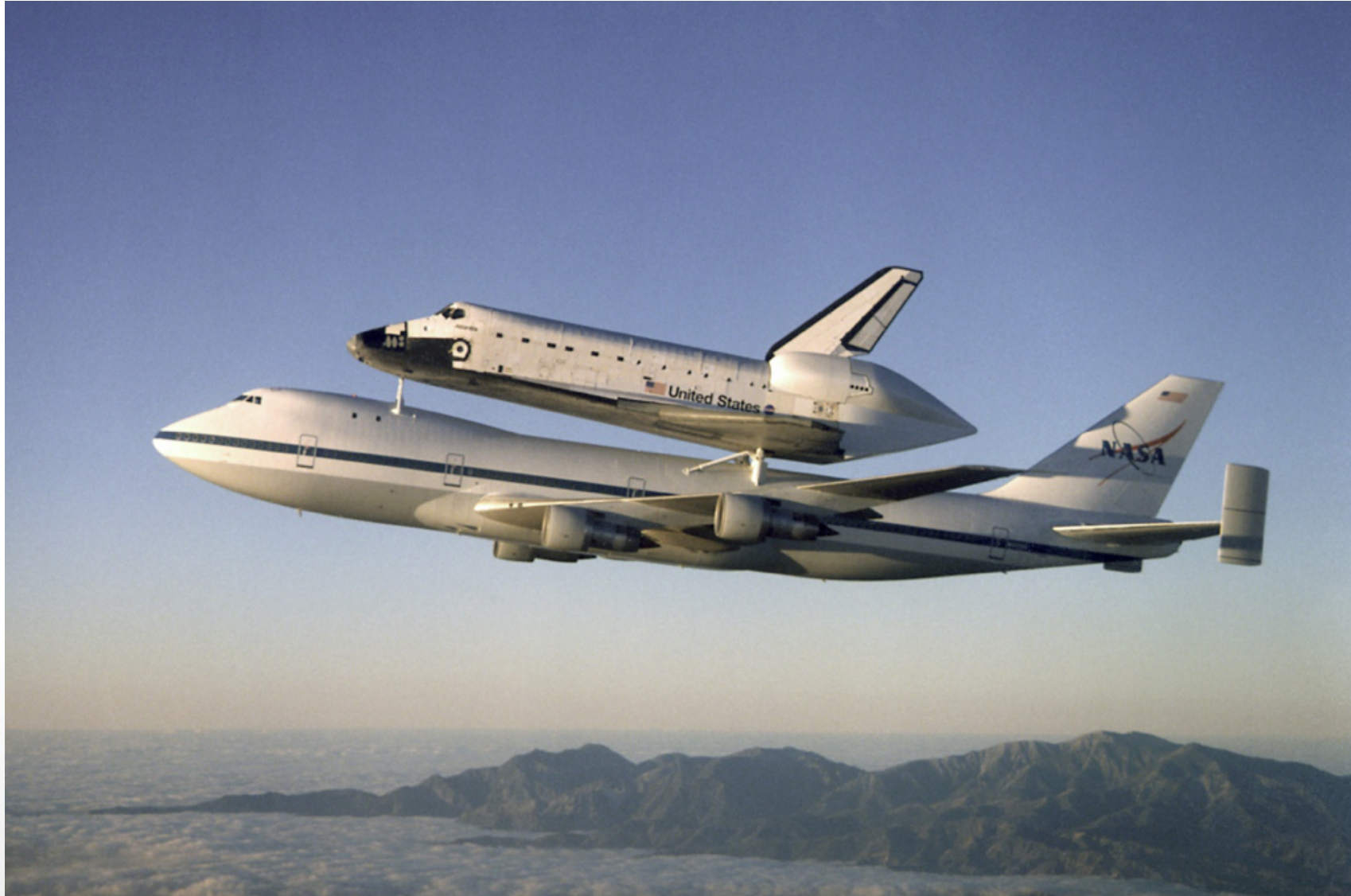


GOAL FOR EXASCALE HPC

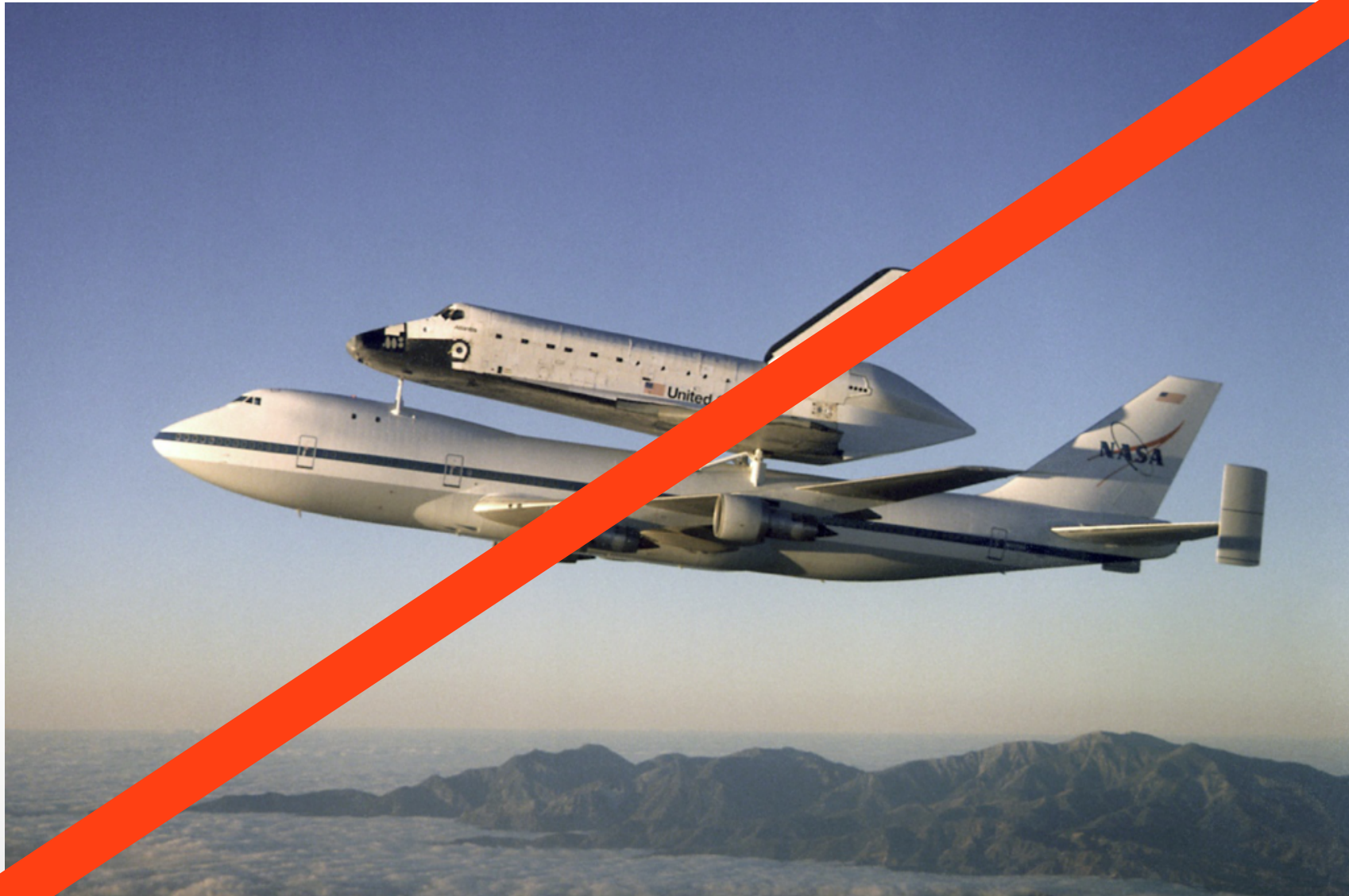


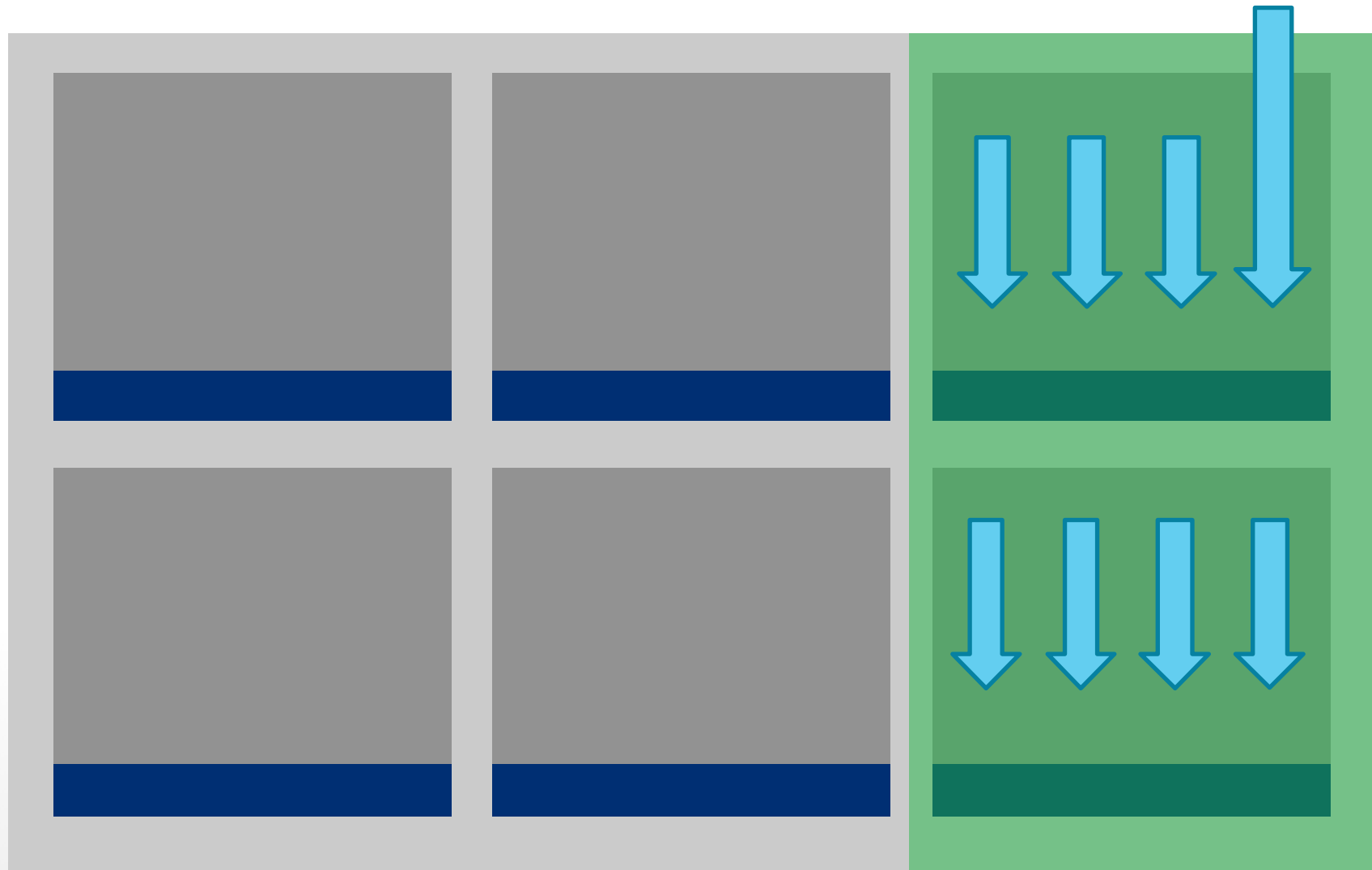


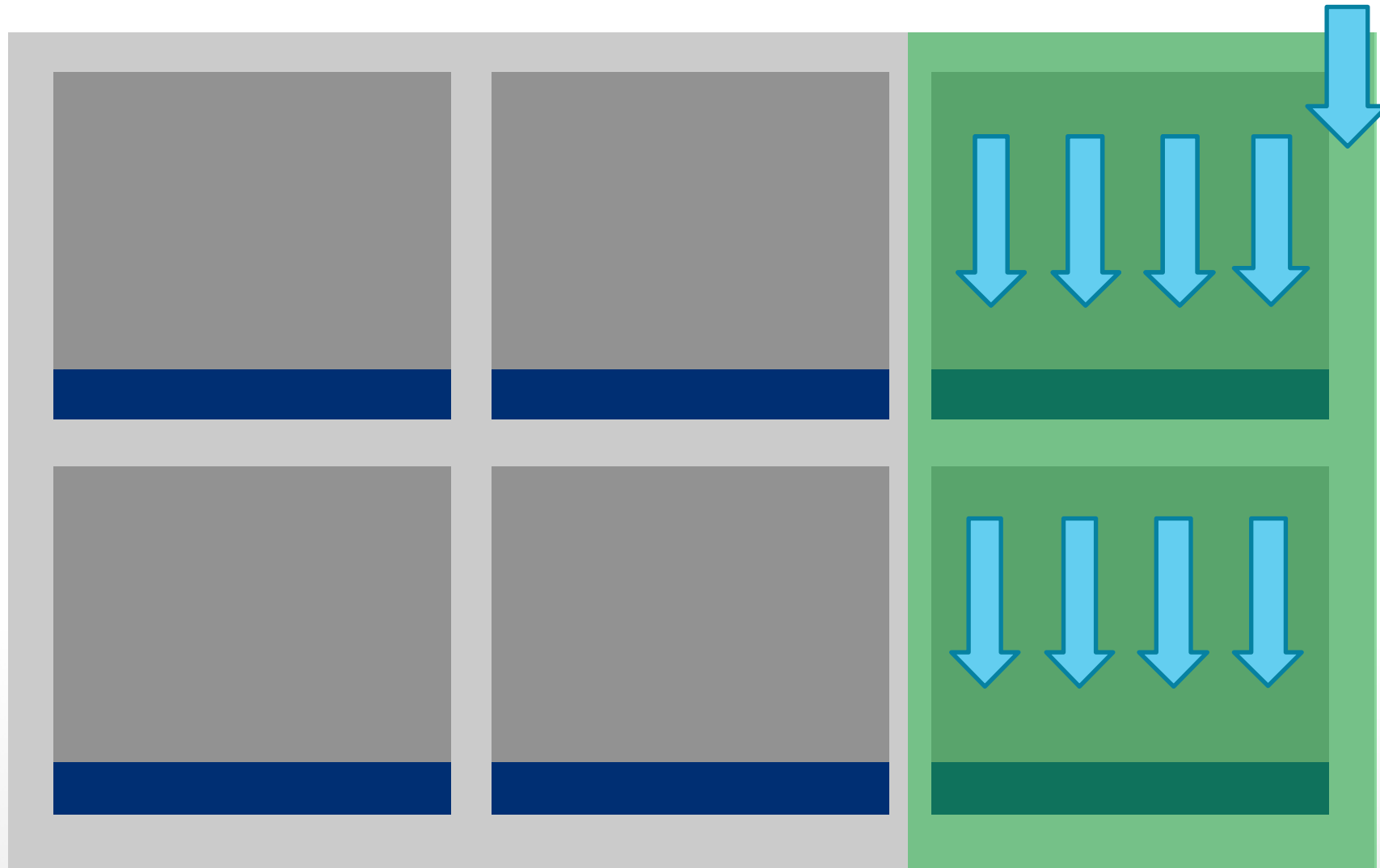
SMALL? PREDICTABLE?

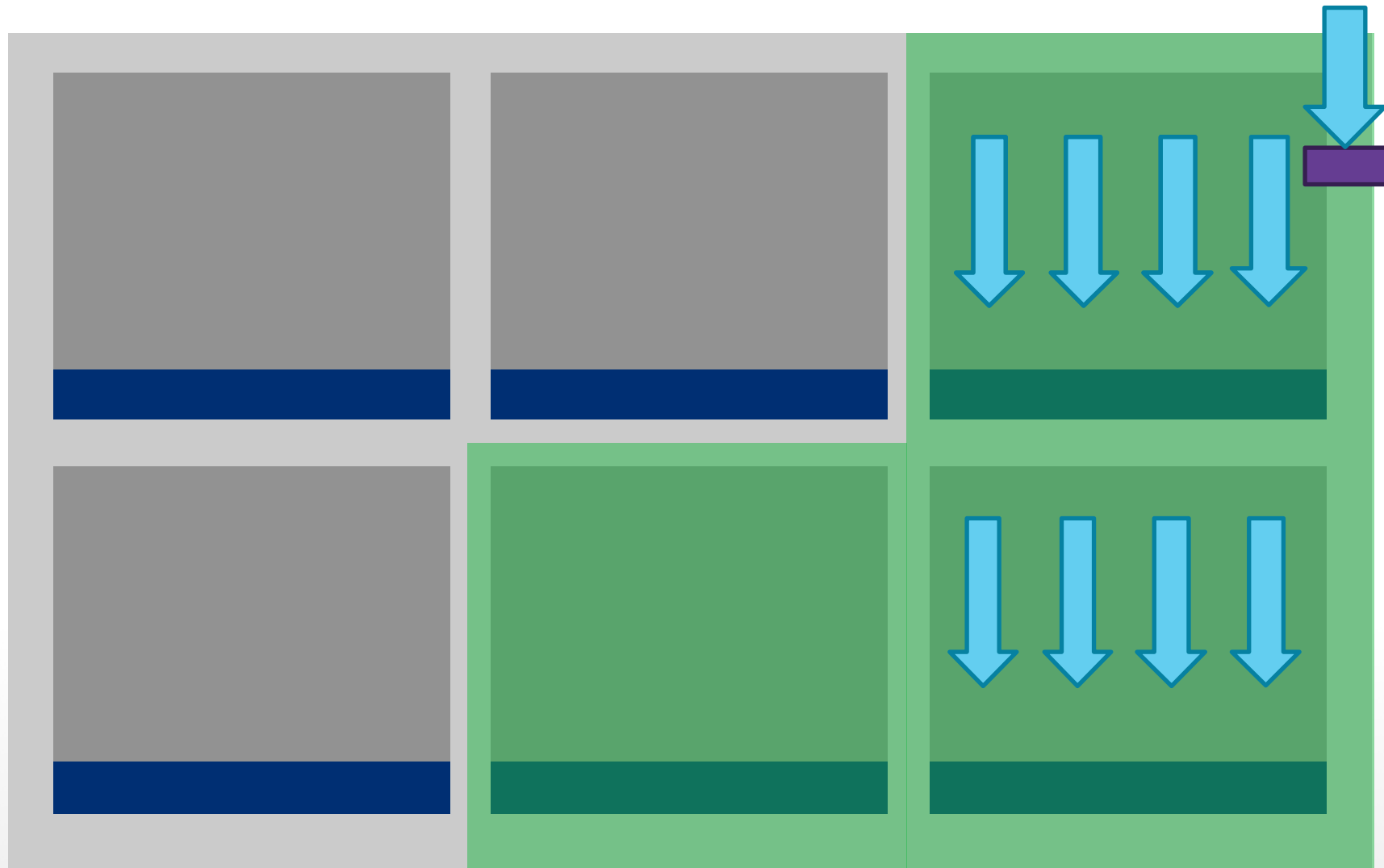


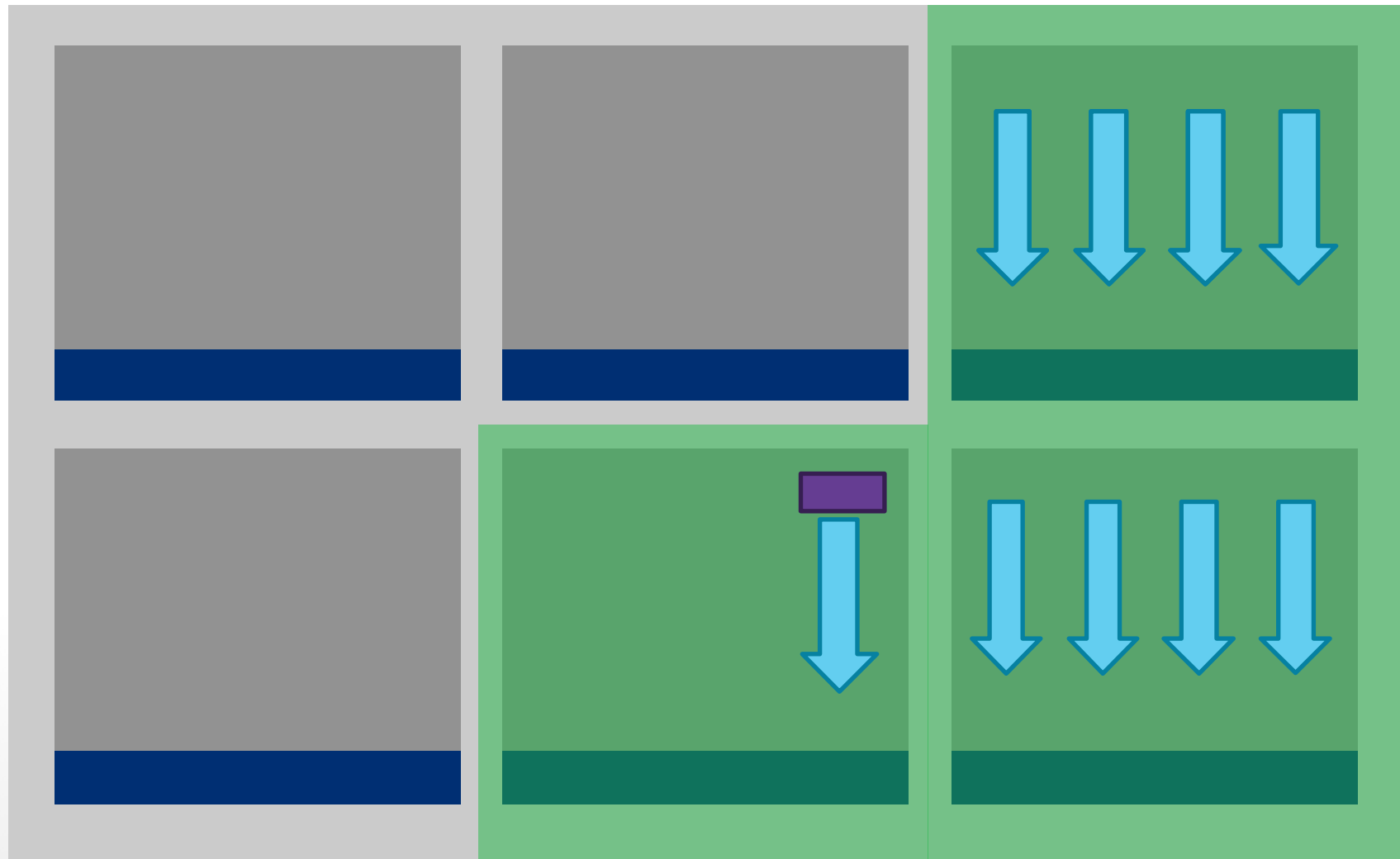
SMALL? PREDICTABLE?



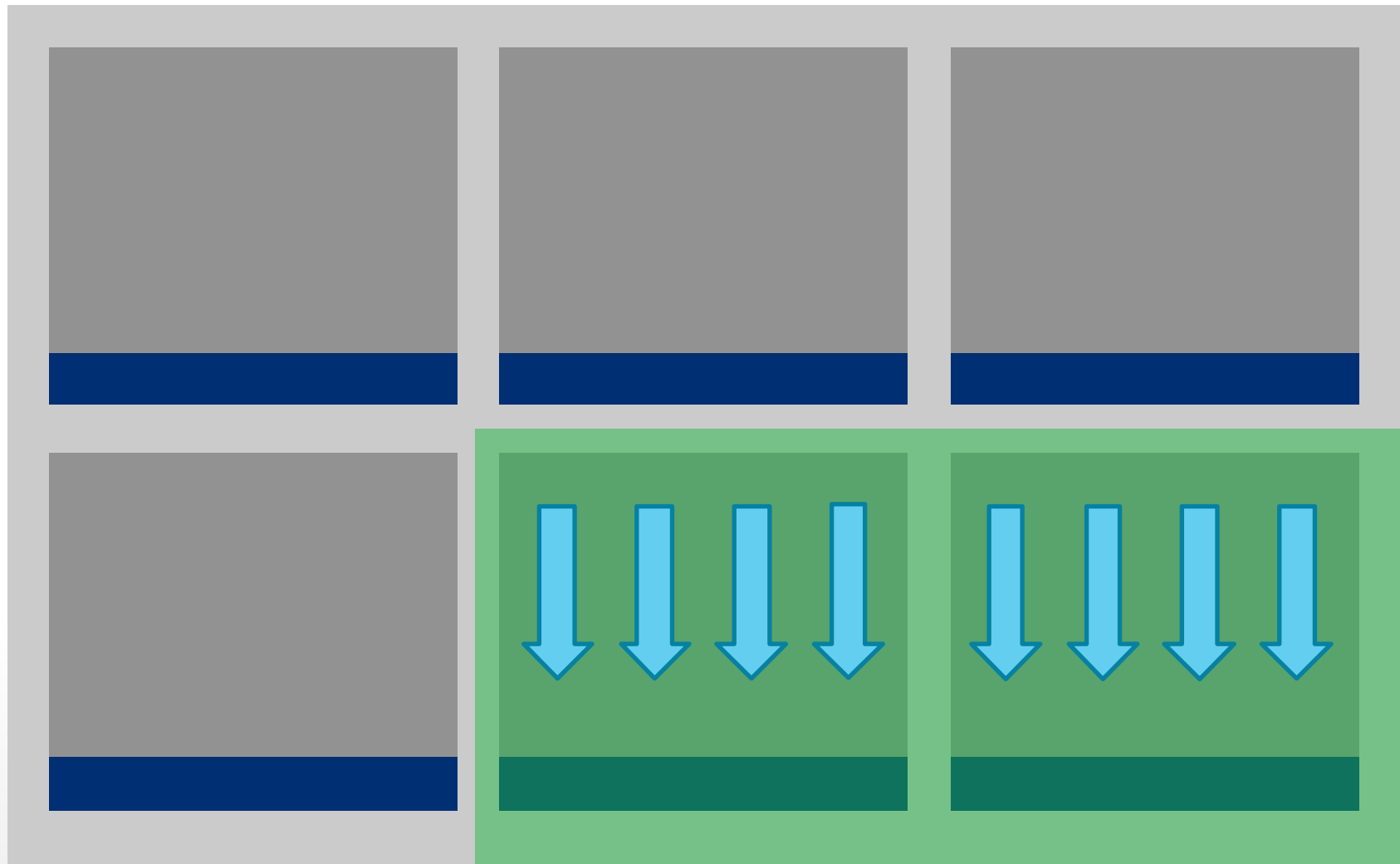




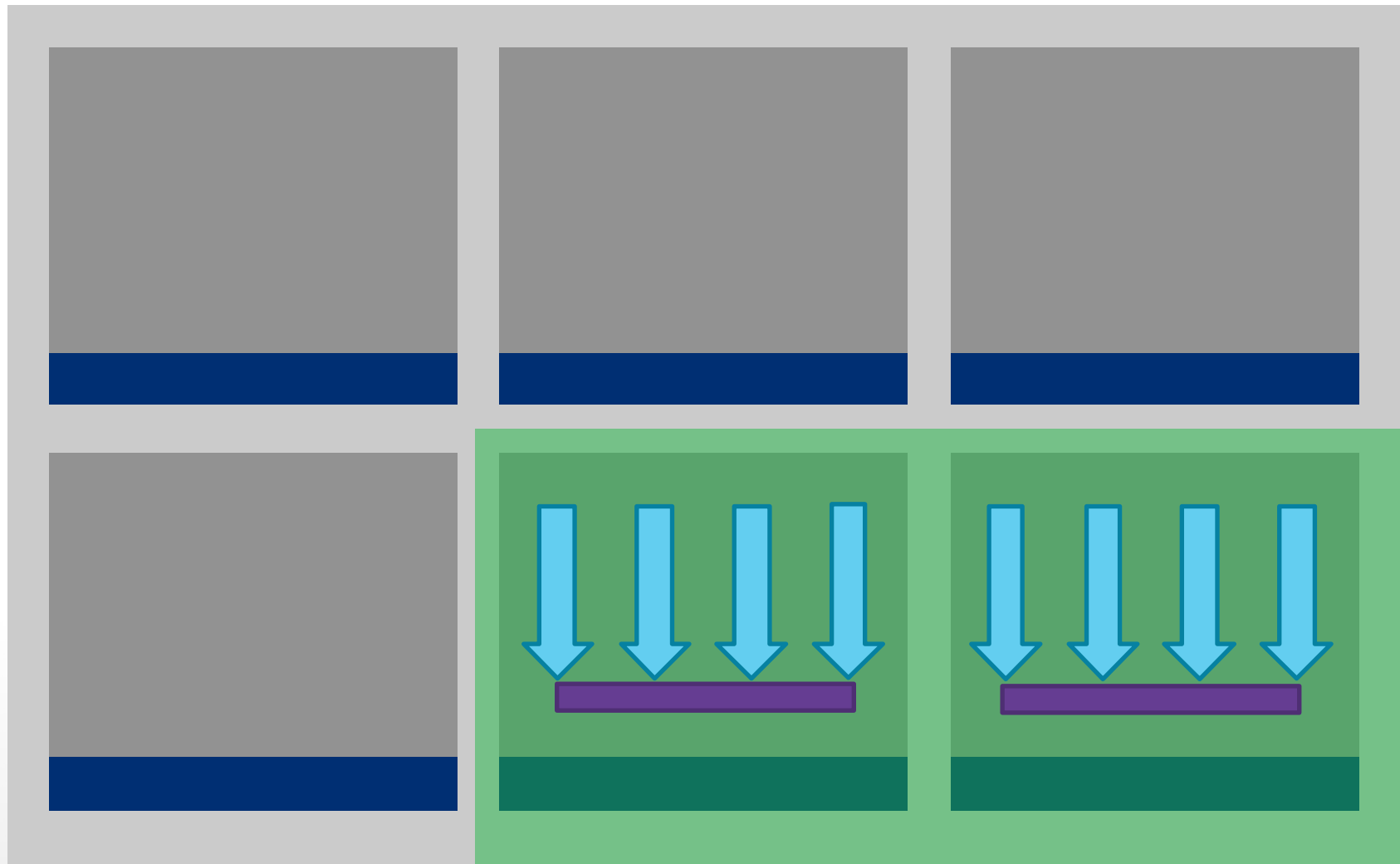




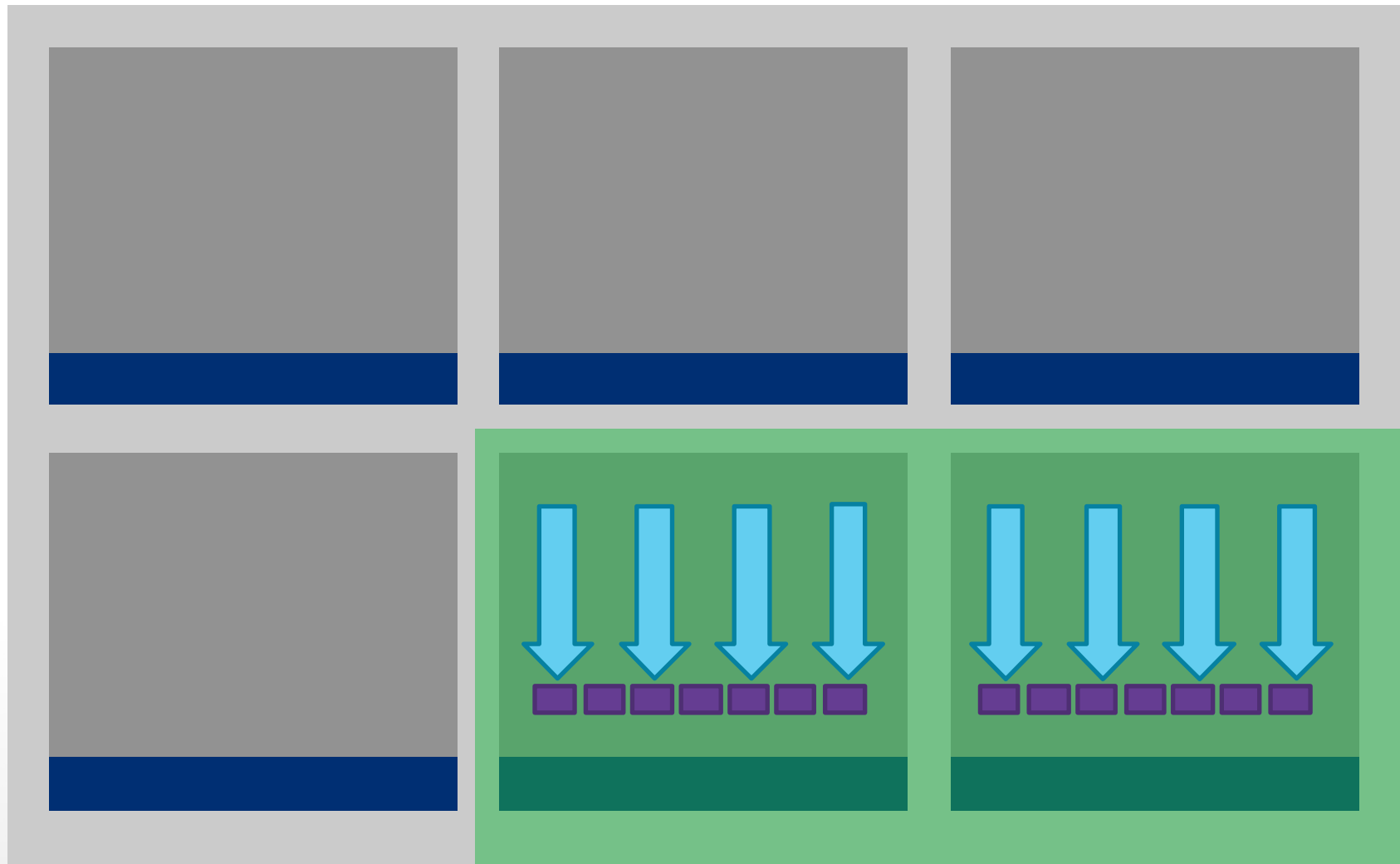
REDUNDANT CHECKPOINT



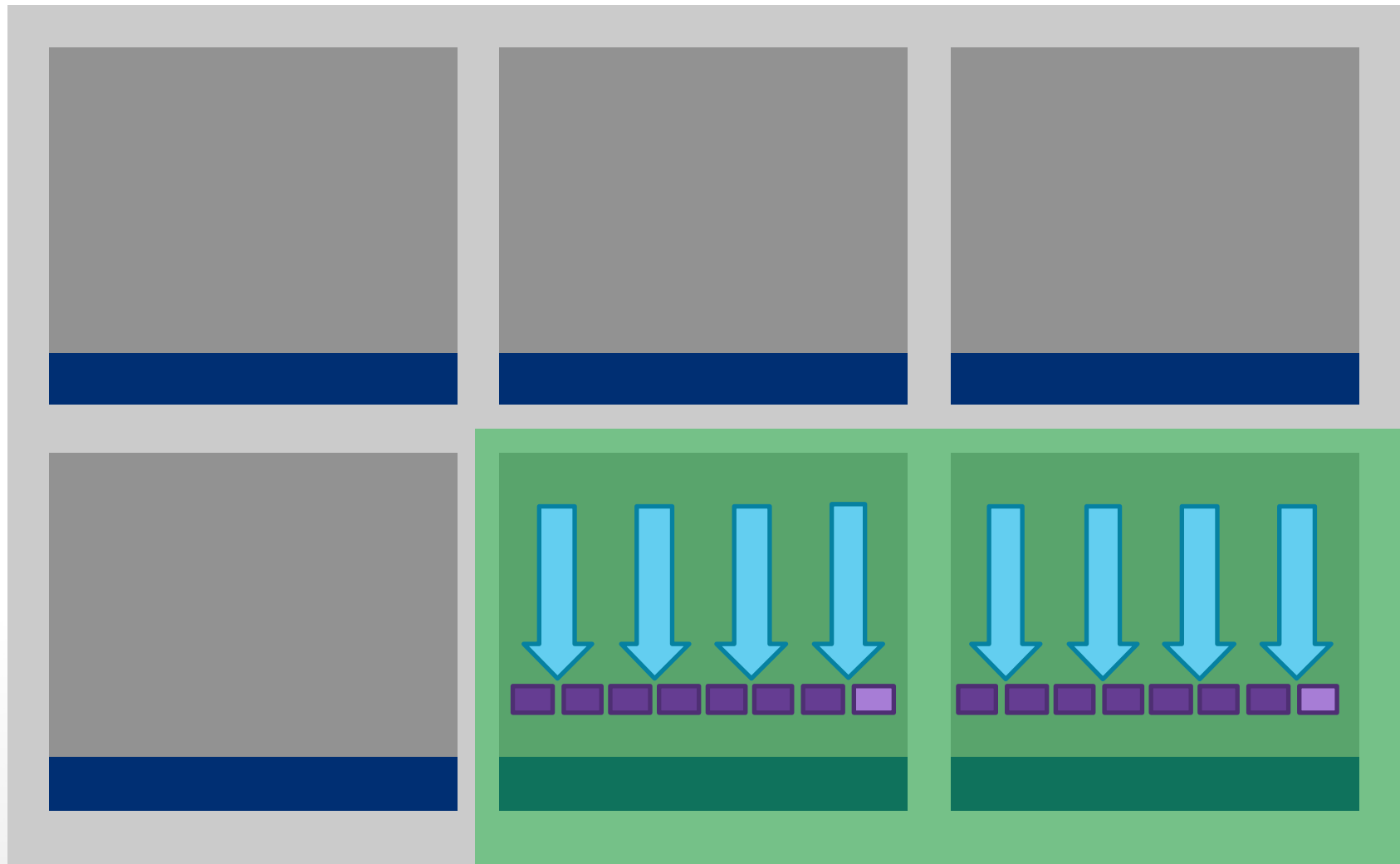
REDUNDANT CHECKPOINT



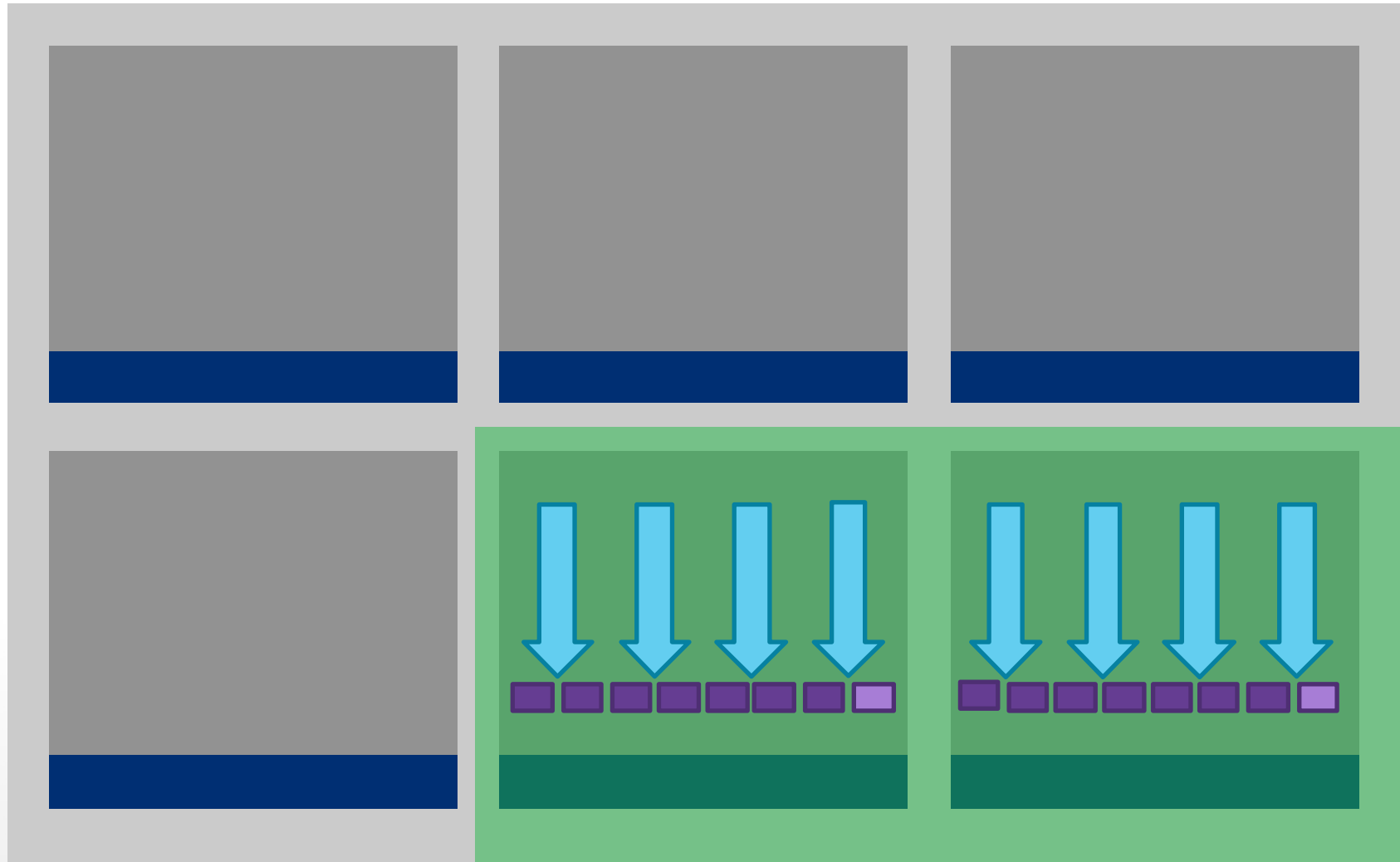
REDUNDANT CHECKPOINT



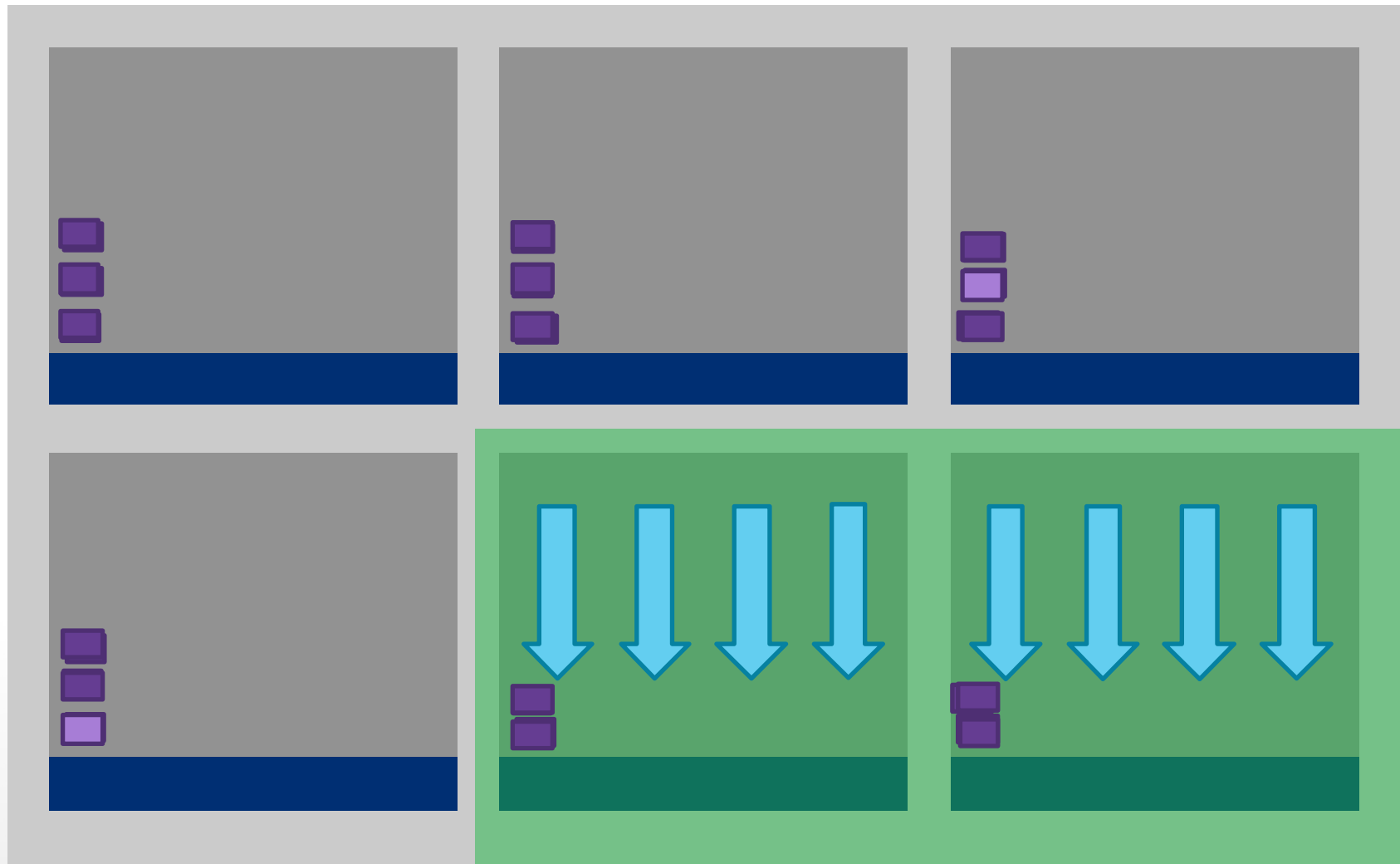
REDUNDANT CHECKPOINT



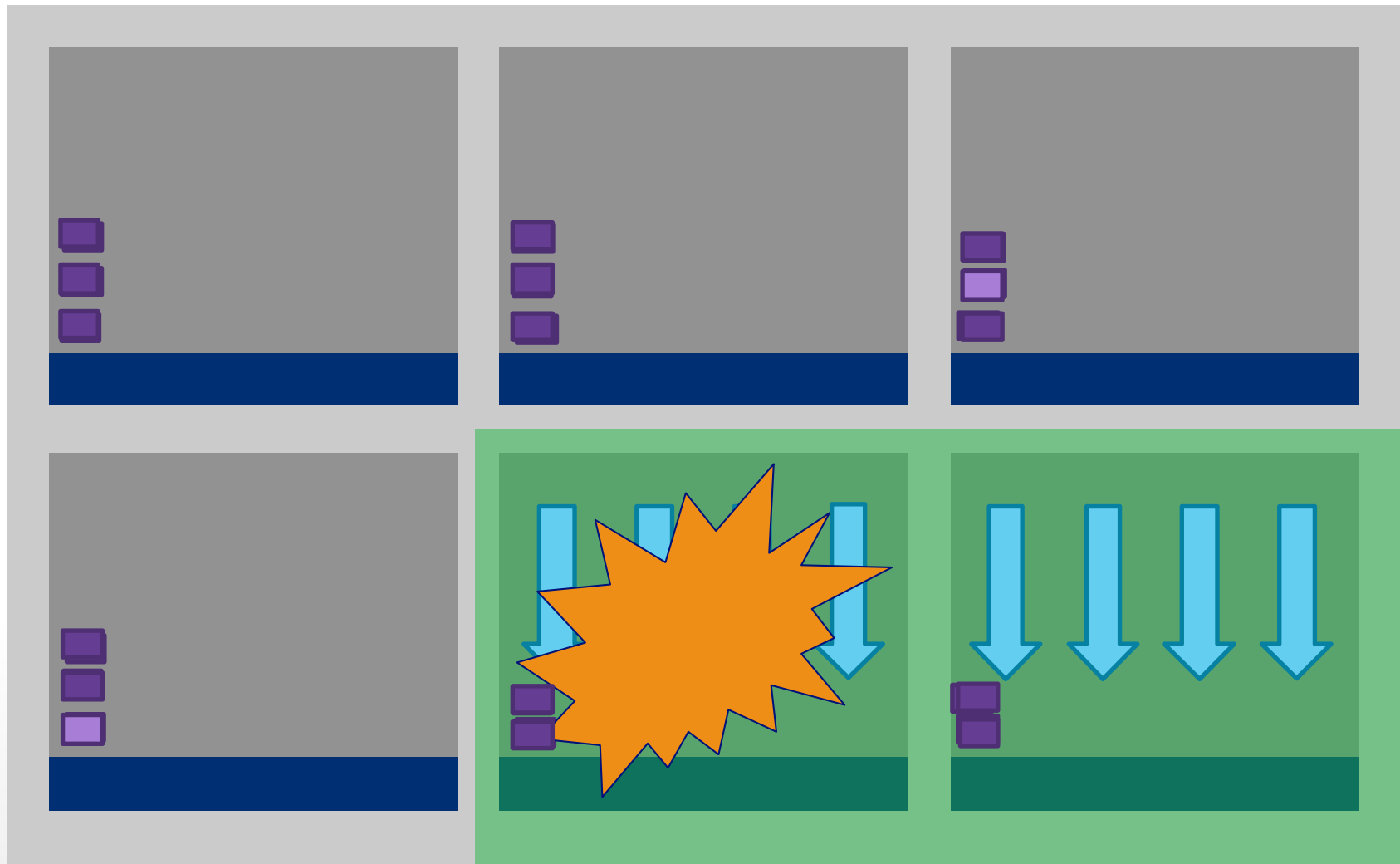
REDUNDANT CHECKPOINT

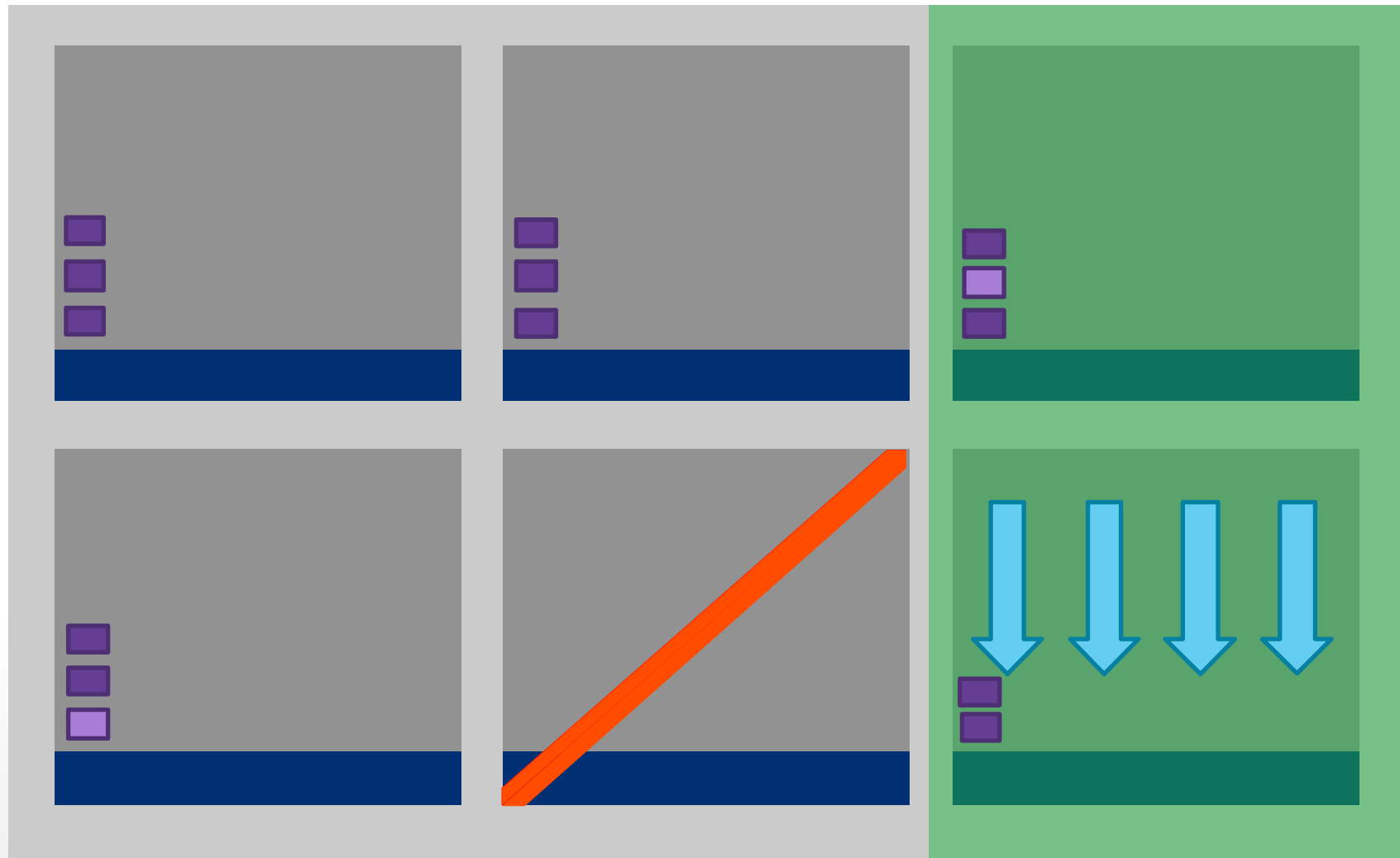


REDUNDANT CHECKPOINT

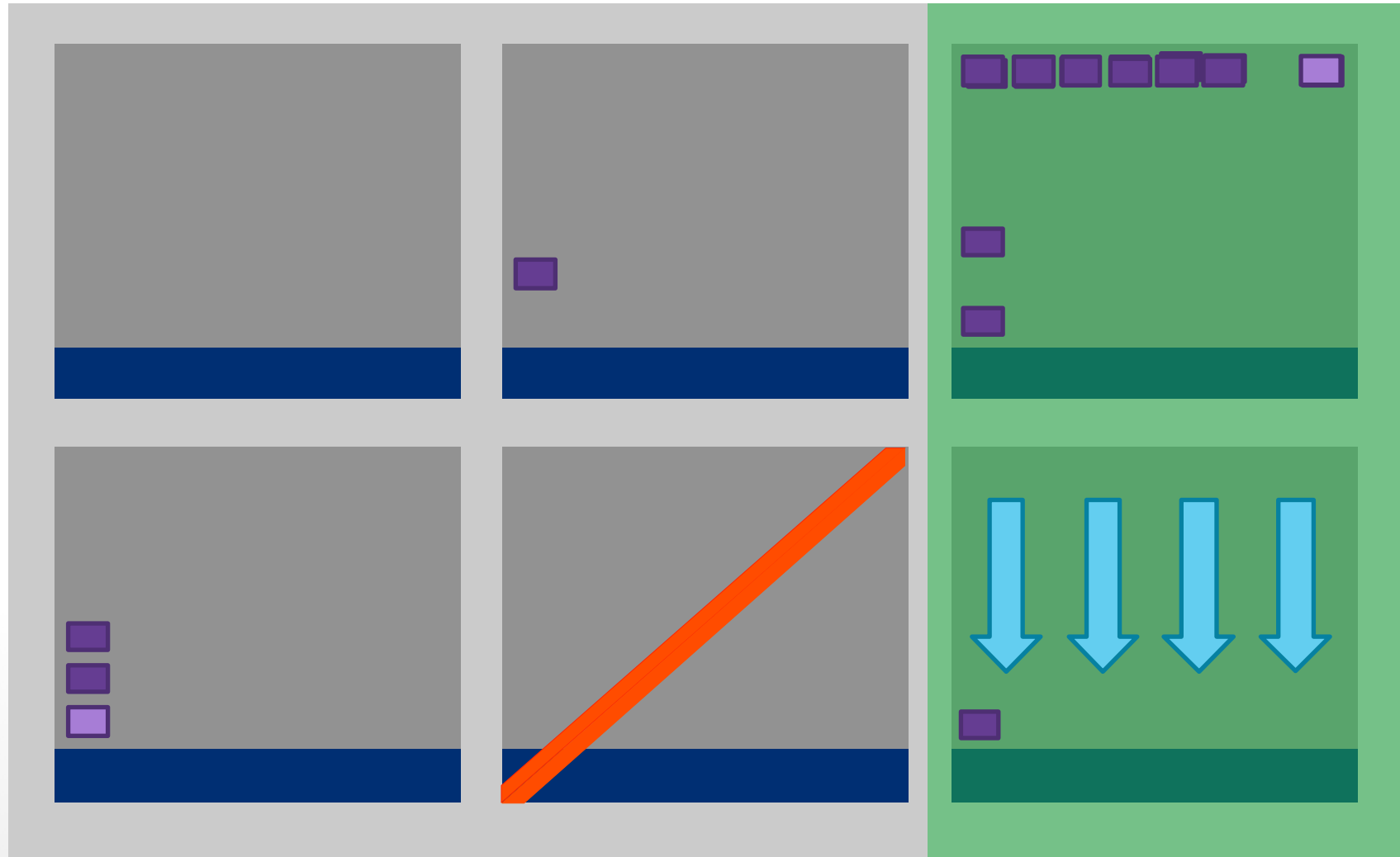


REDUNDANT CHECKPOINT

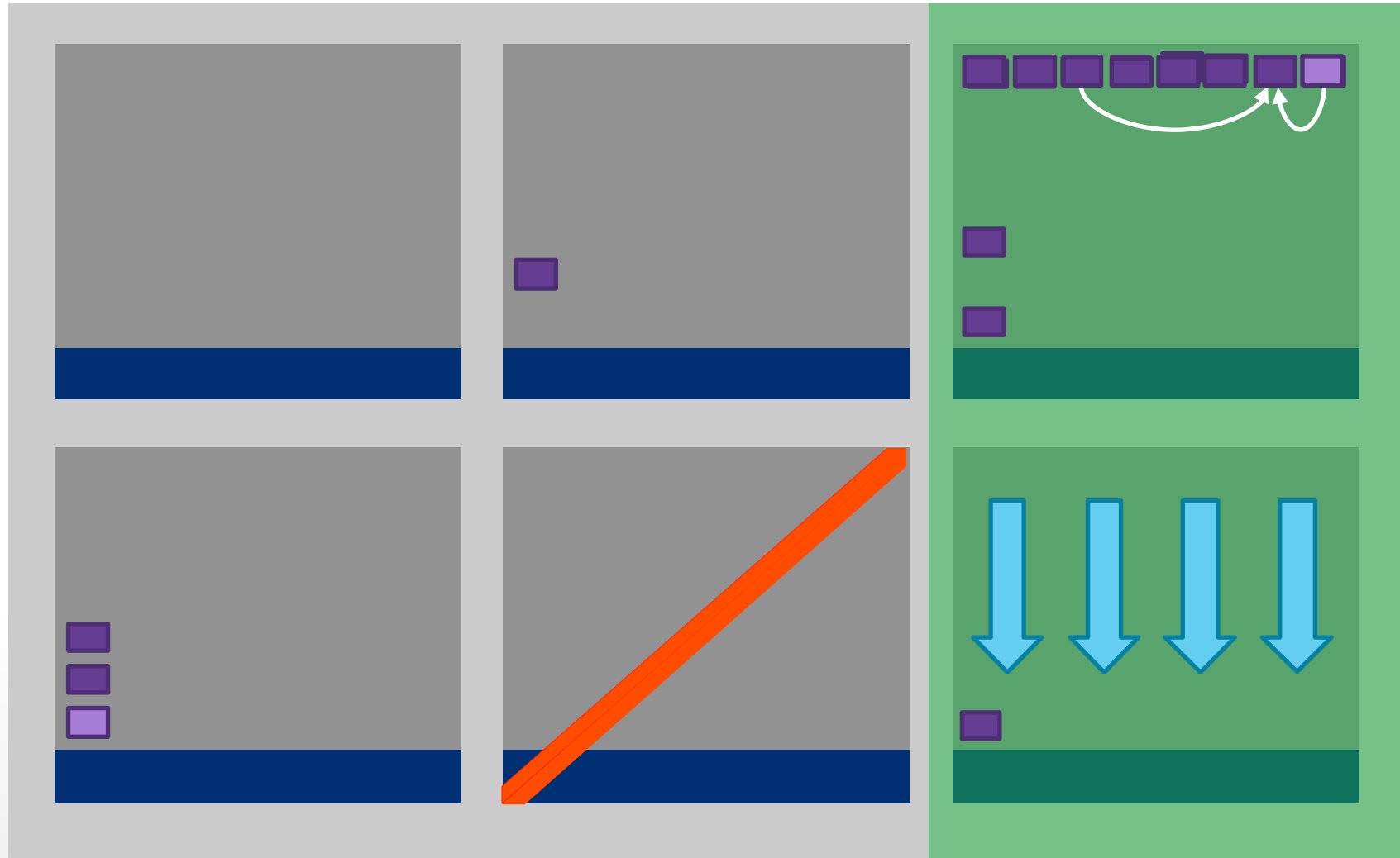




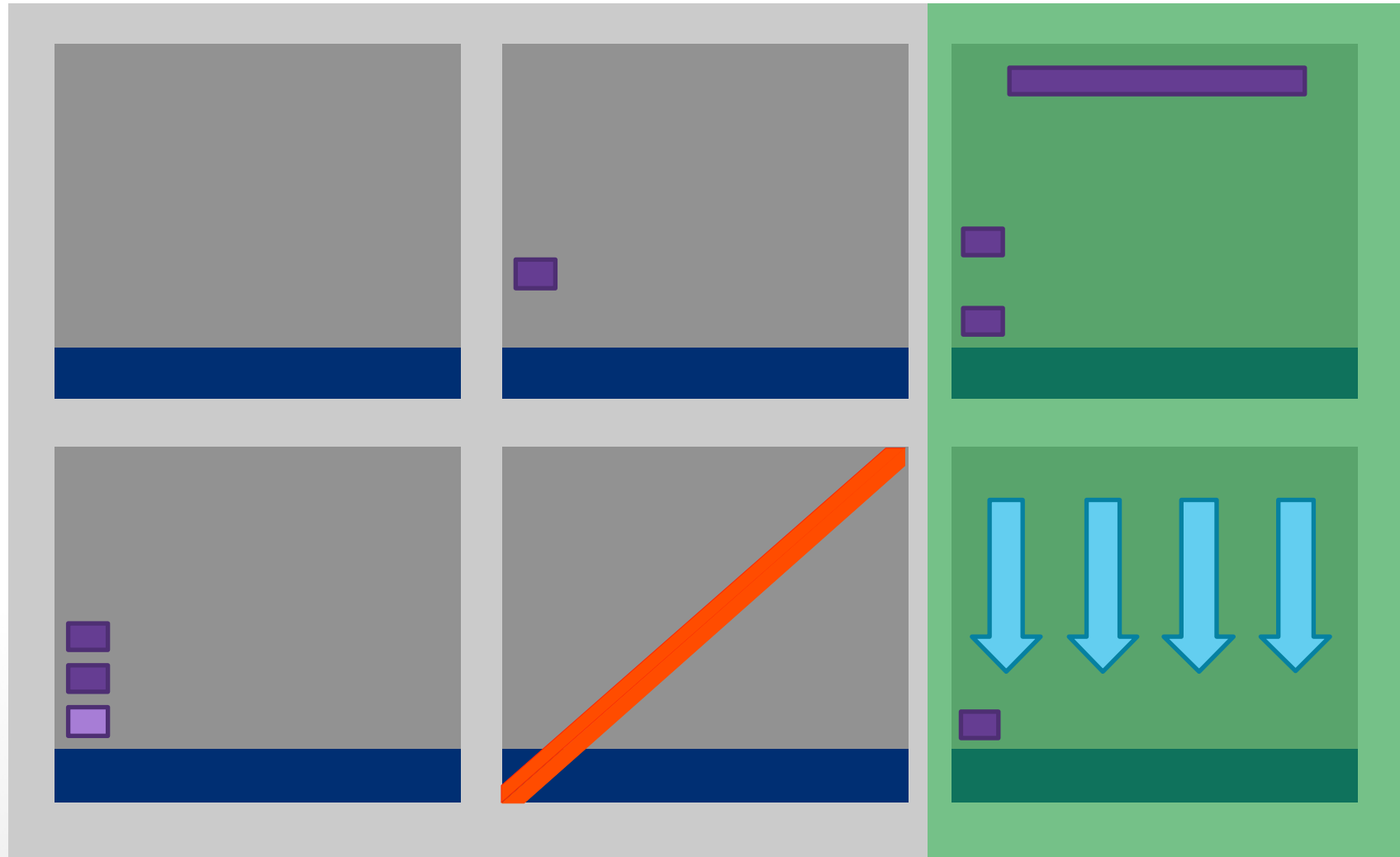
REDUNDANT CHECKPOINT



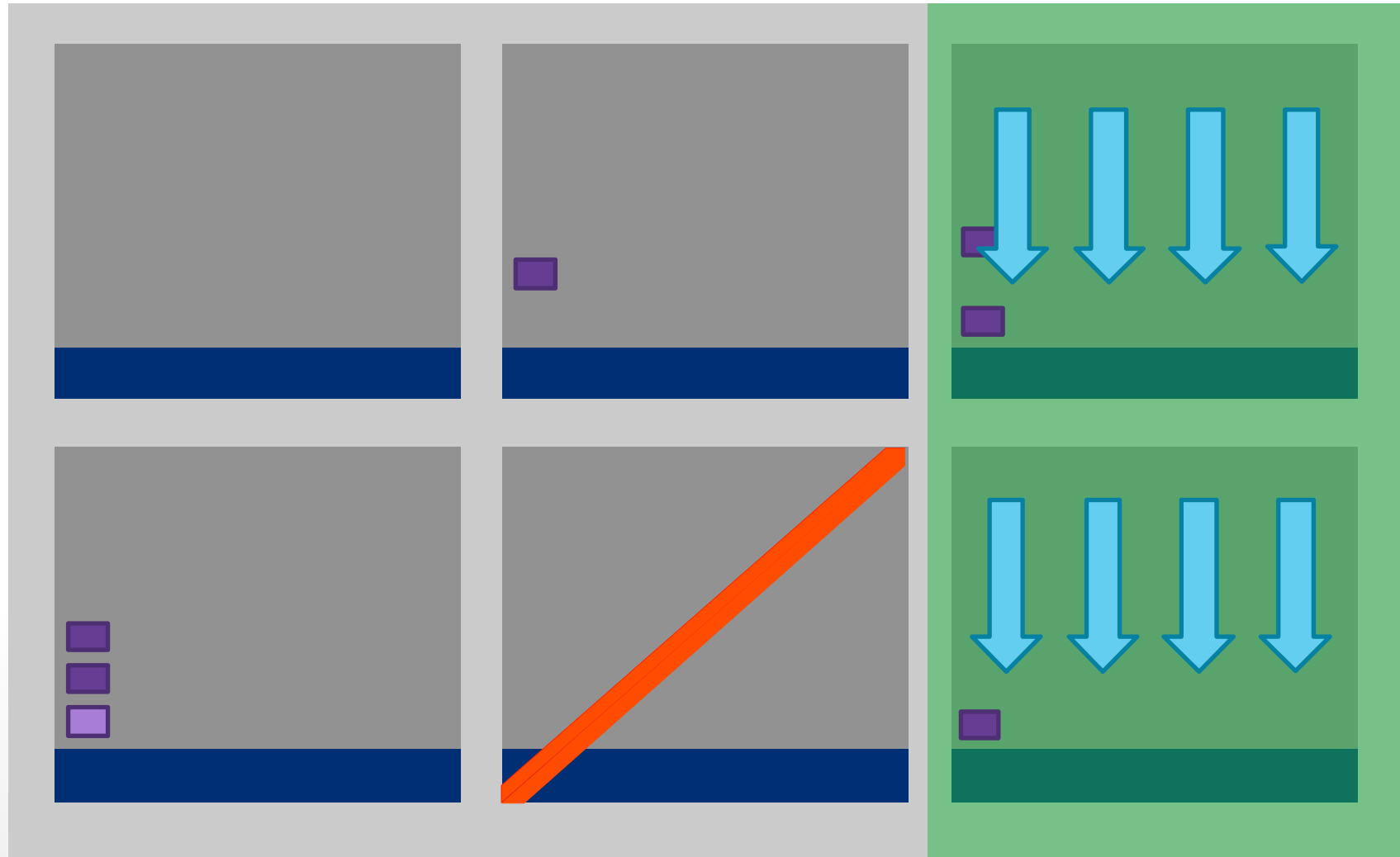
REDUNDANT CHECKPOINT



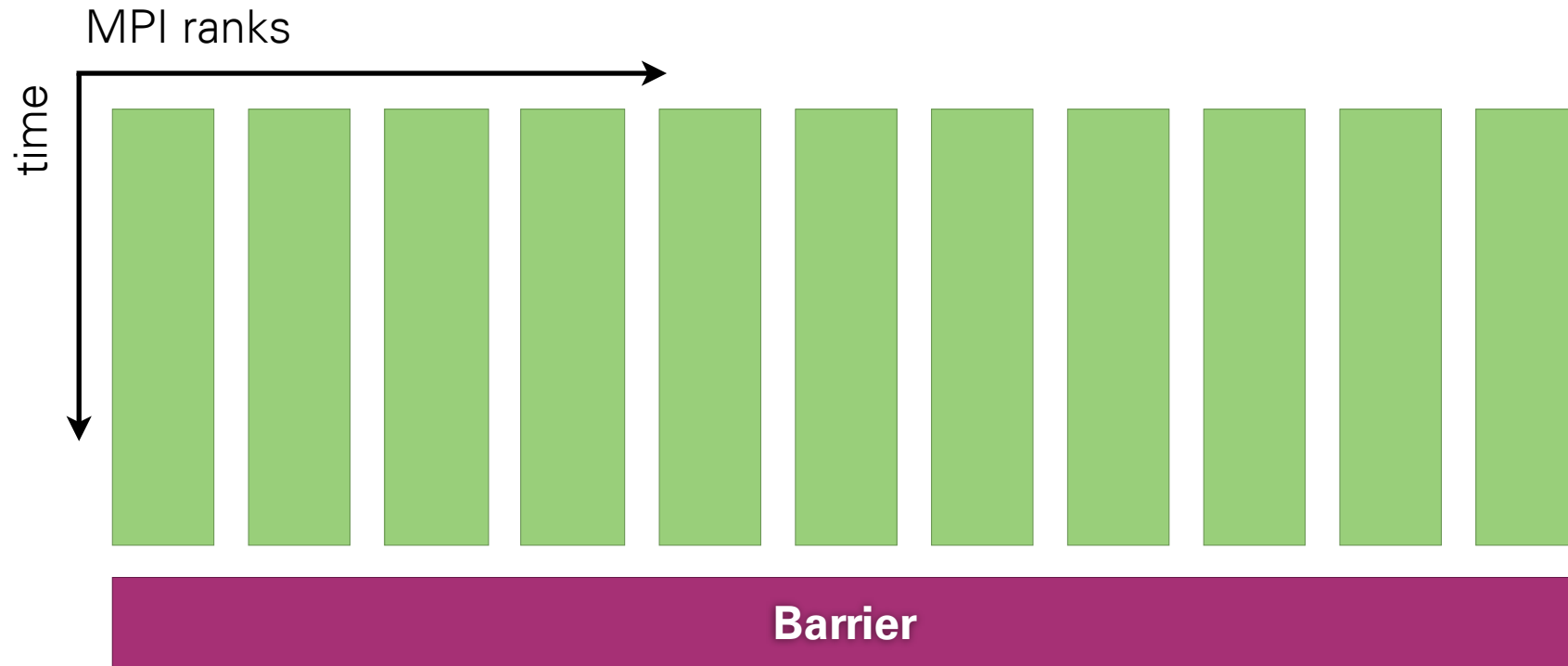
REDUNDANT CHECKPOINT

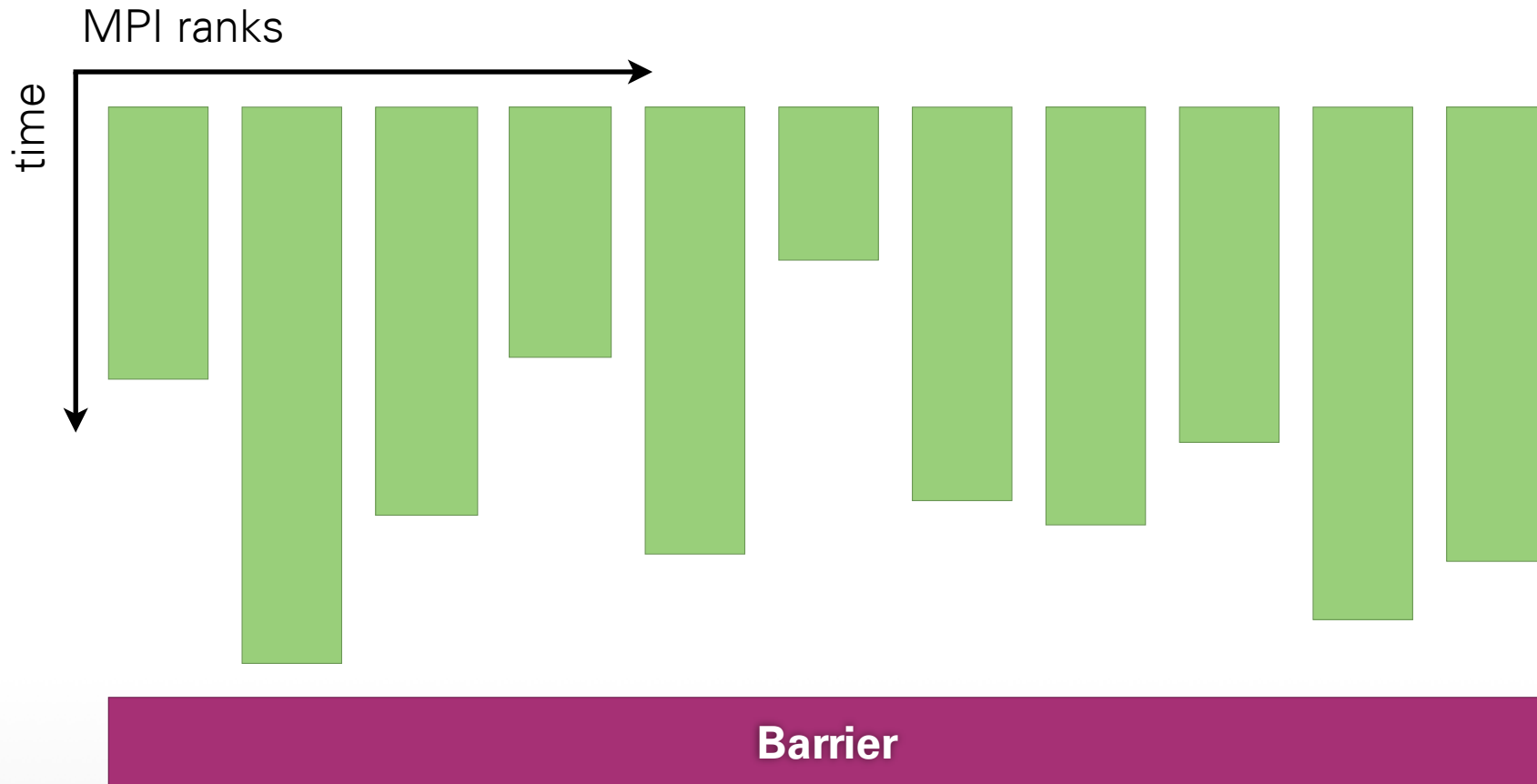


REDUNDANT CHECKPOINT

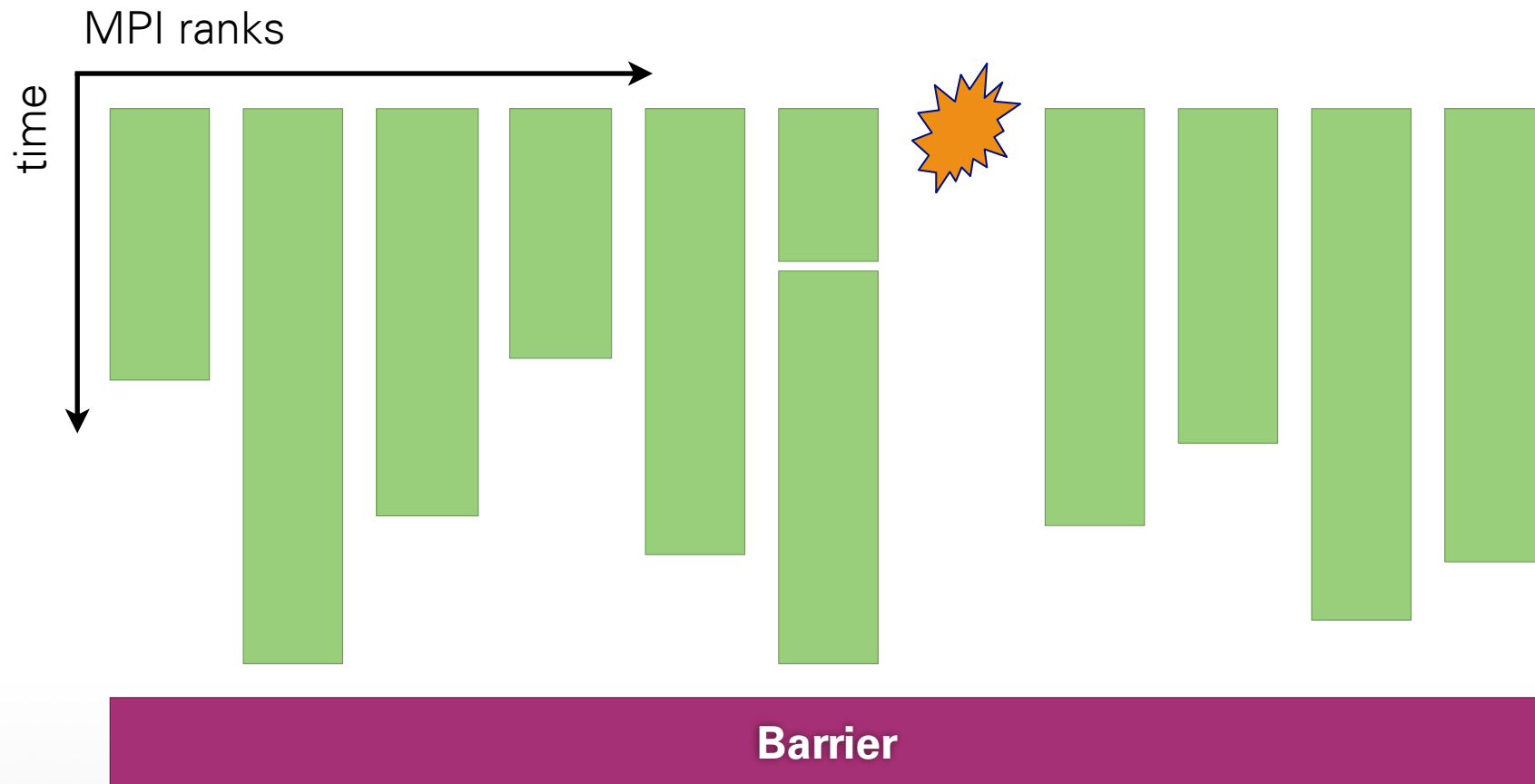


EXPERIMENTS: IMBALANCES, OVERDECOMPOSITION AND OVERSUBSCRIPTION



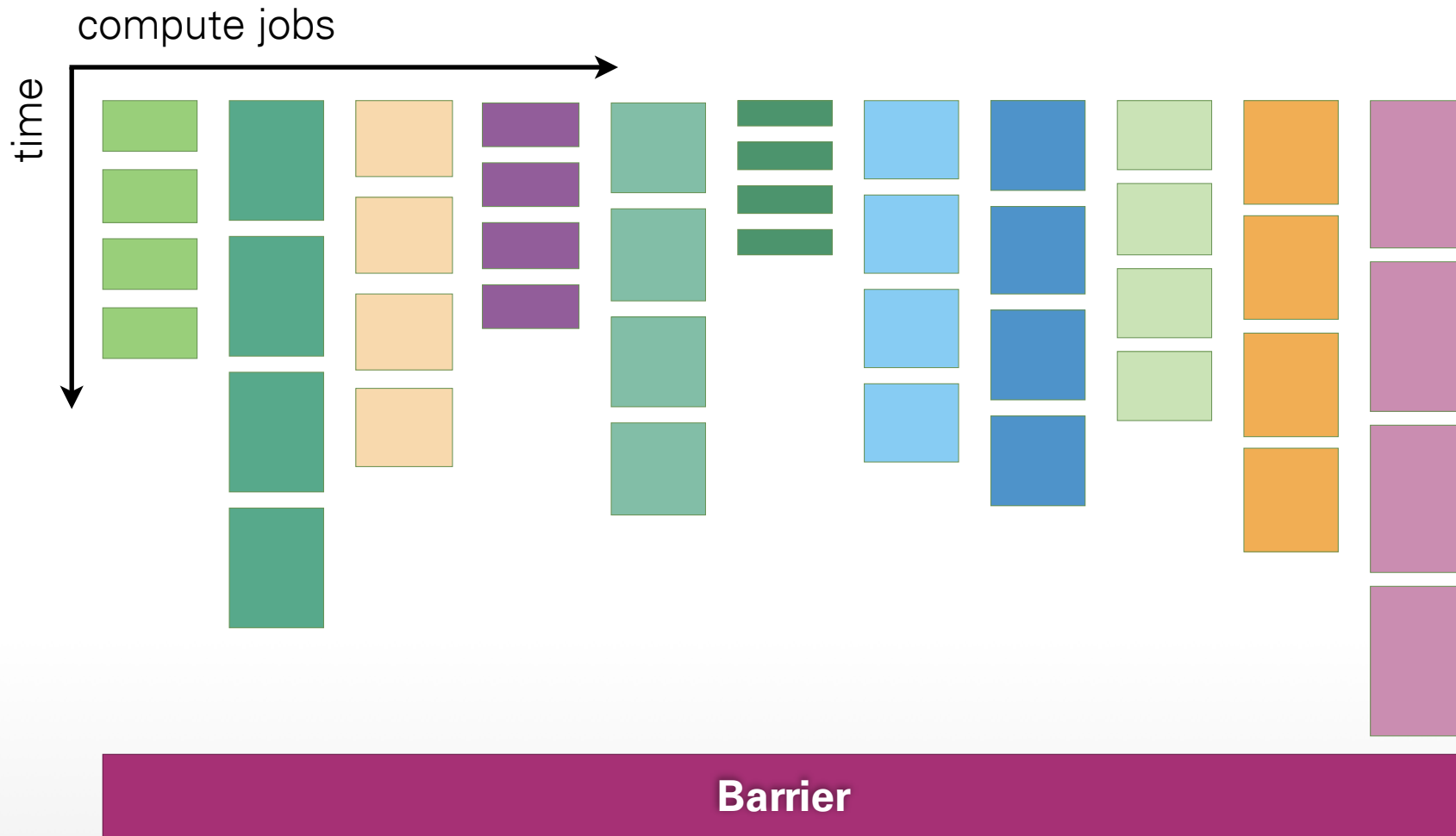


Imbalance in application workload



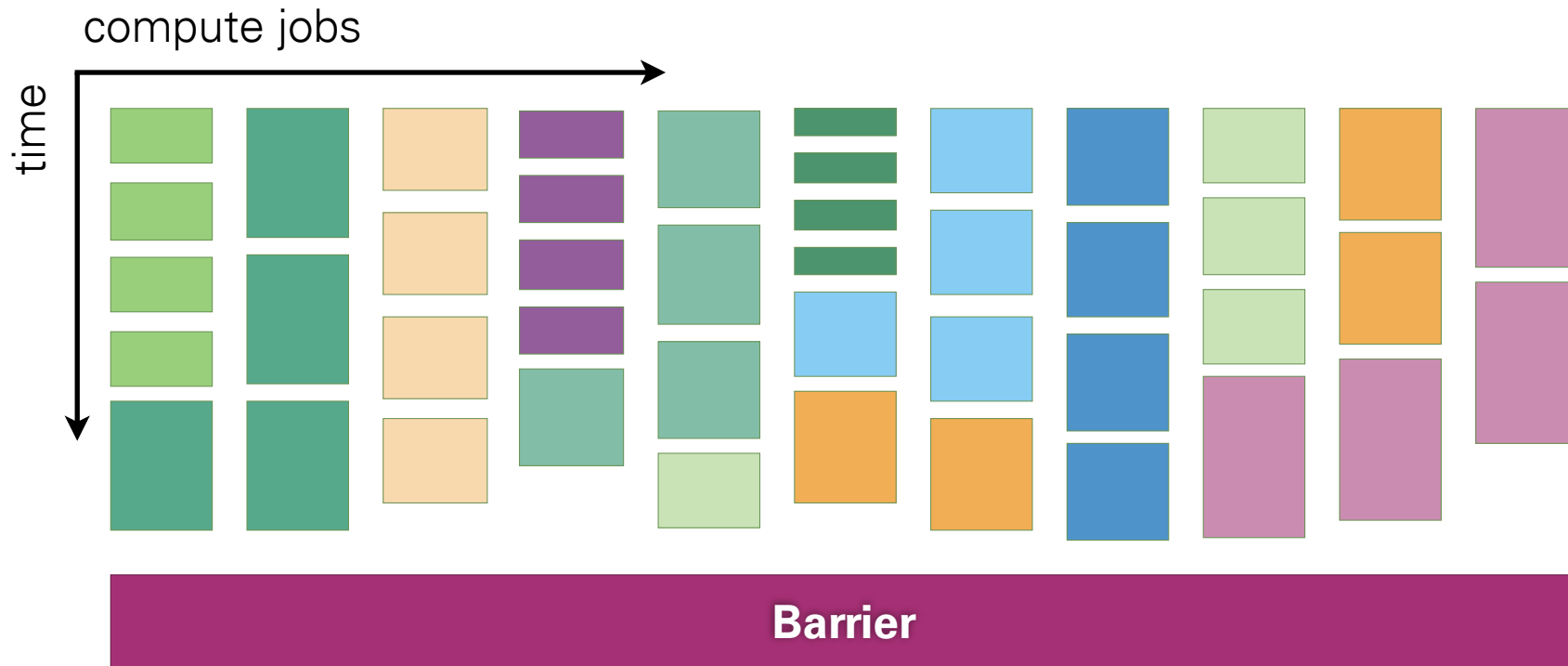
Reassign work to react to node failure

SPLITTING BIG JOBS



overdecomposition & "oversubscription"

SMALL JOBS (NO DEPS)



Execute small jobs in parallel (if possible)



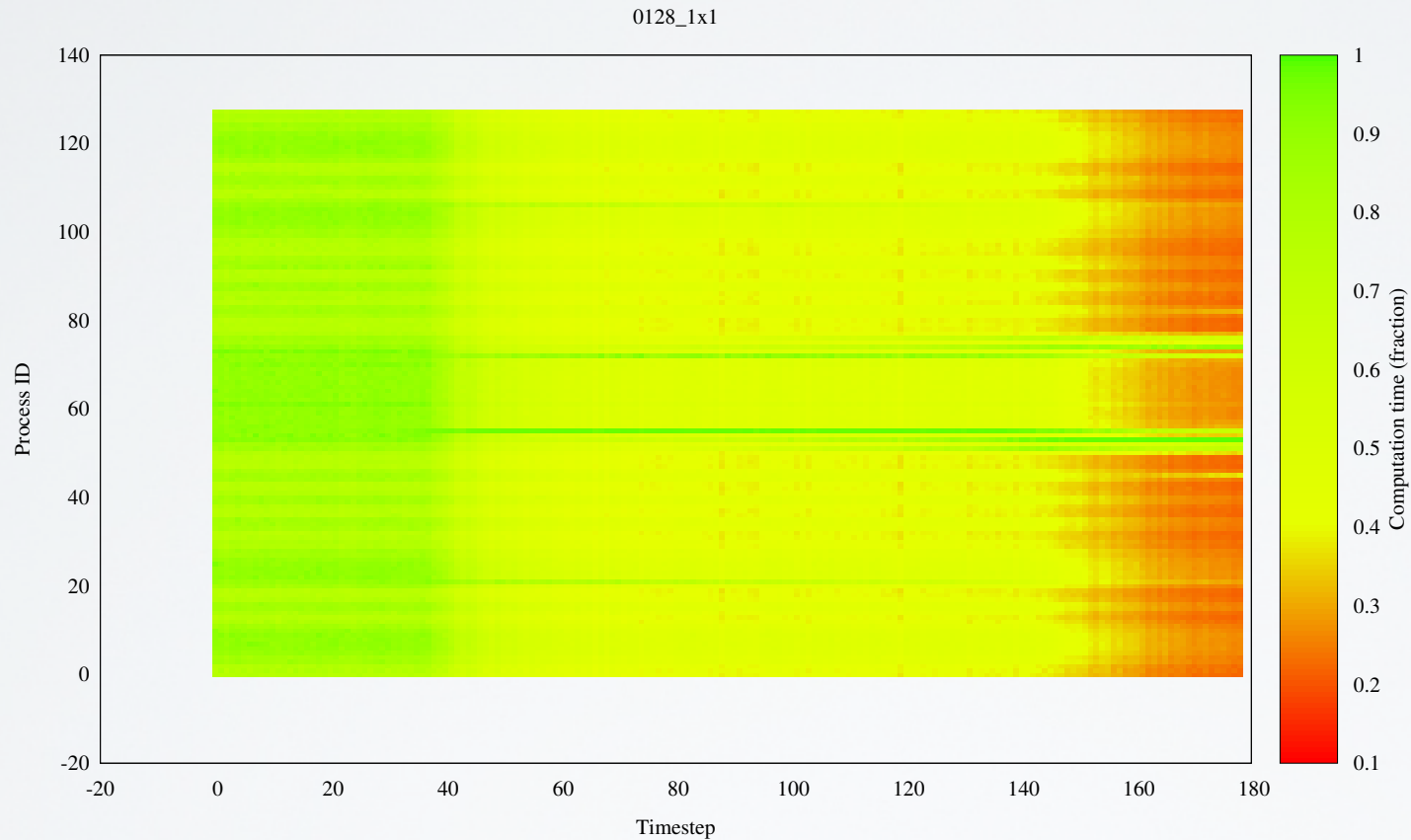
Unbalanced compute times of ranks per time step

Application: COSMO-SPECS+FD4



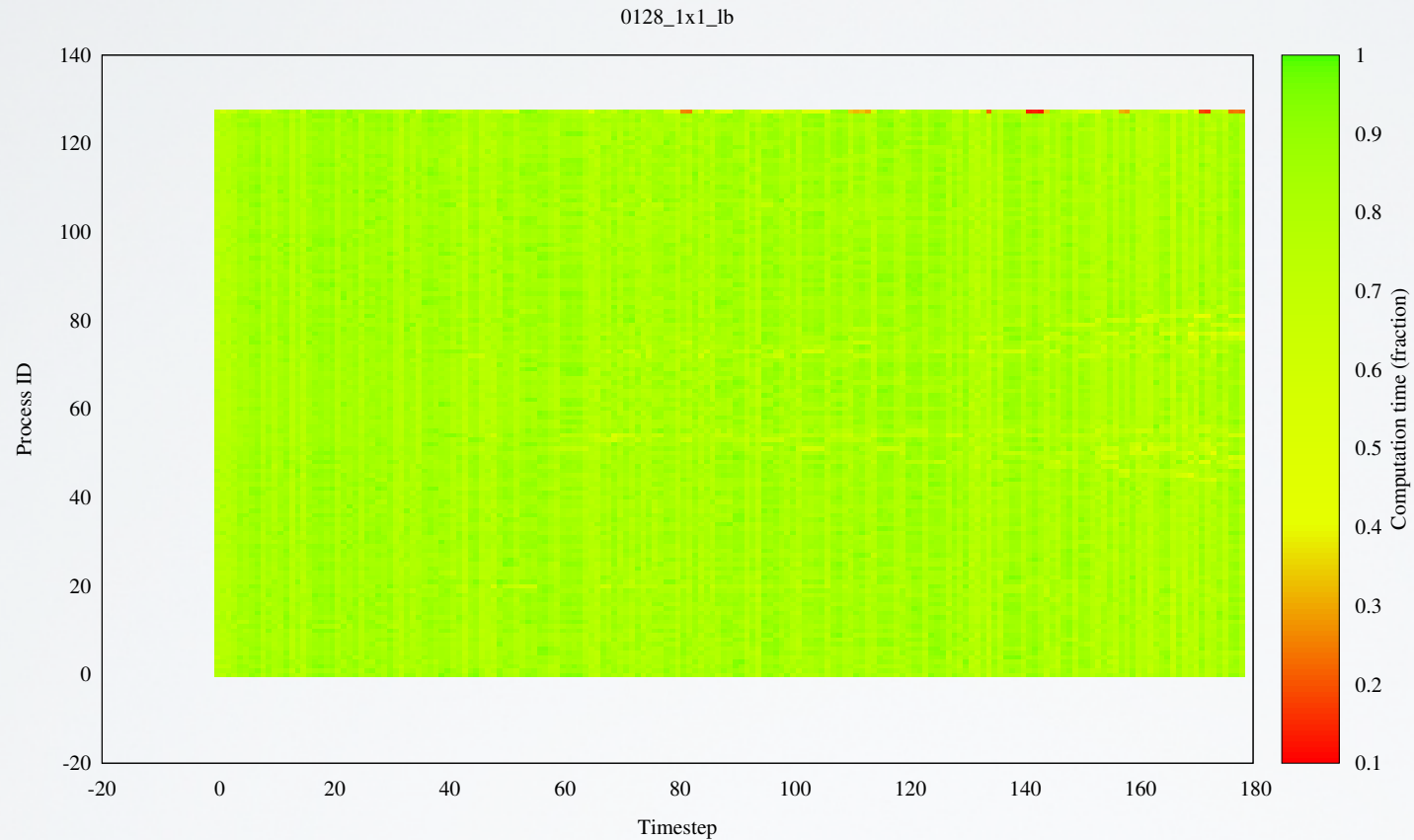
Balanced compute times of ranks per time step

Application: COSMO-SPECS+FD4



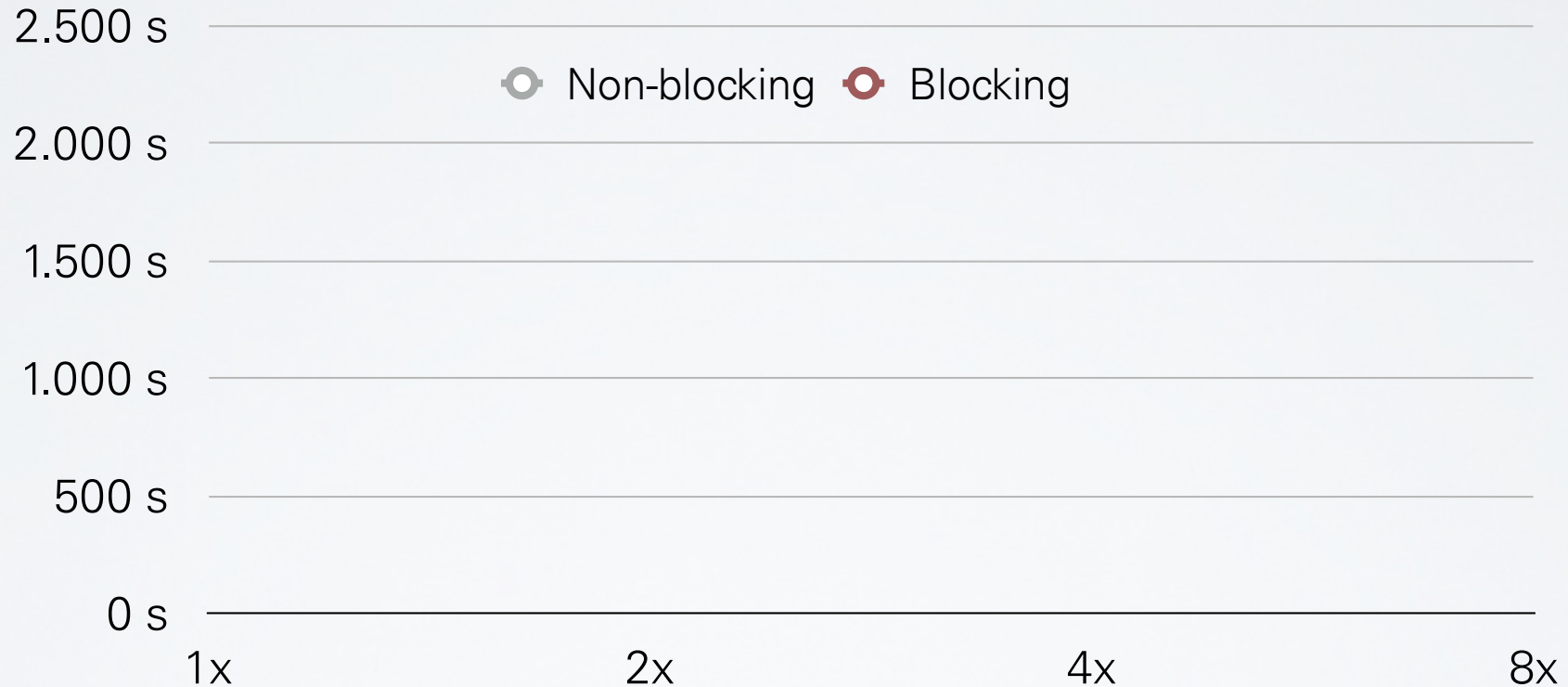
Unbalanced compute times of ranks per time step

Application: COSMO-SPECS+FD4



Balanced compute times of ranks per time step

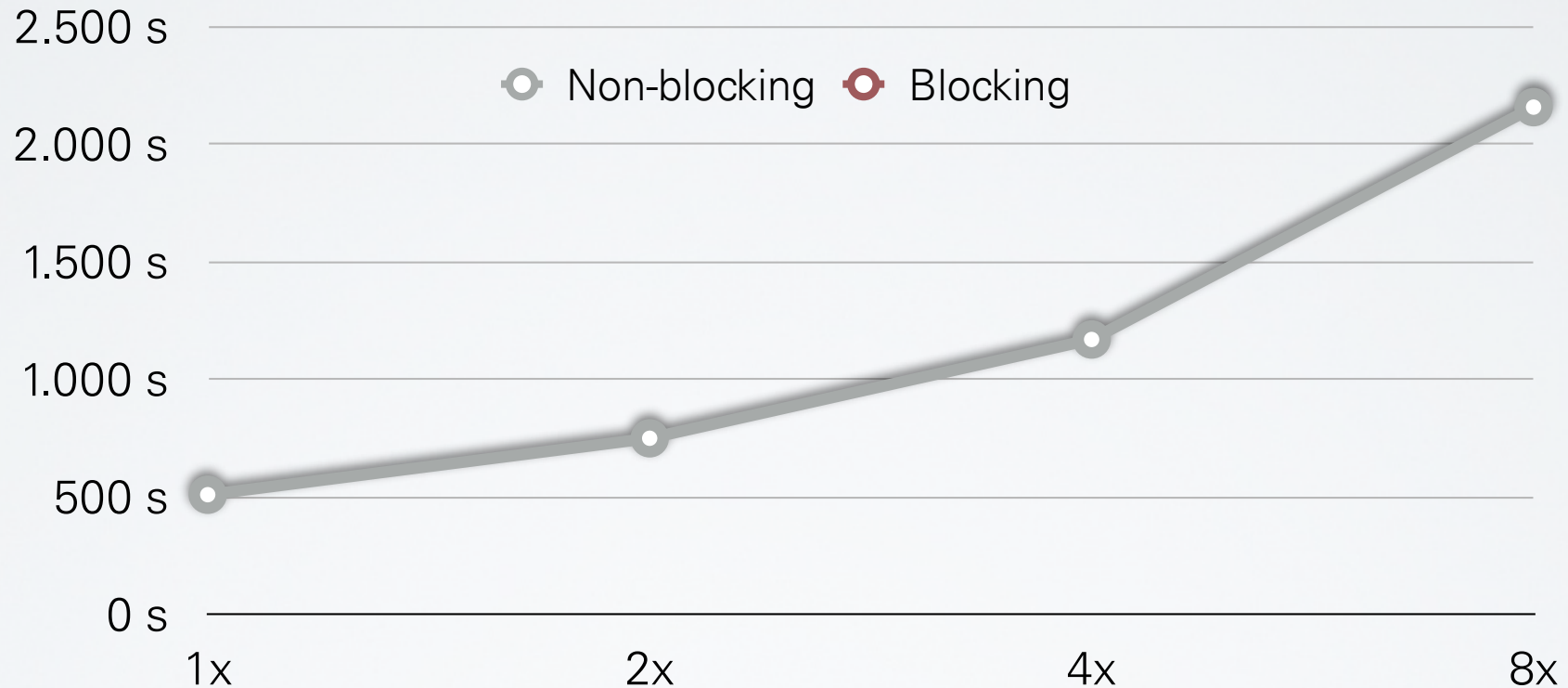
Application: COSMO-SPECS+FD4



Oversubscription factor (more ranks)

Application: COSMO-SPECS+FD4 (no load balancing)

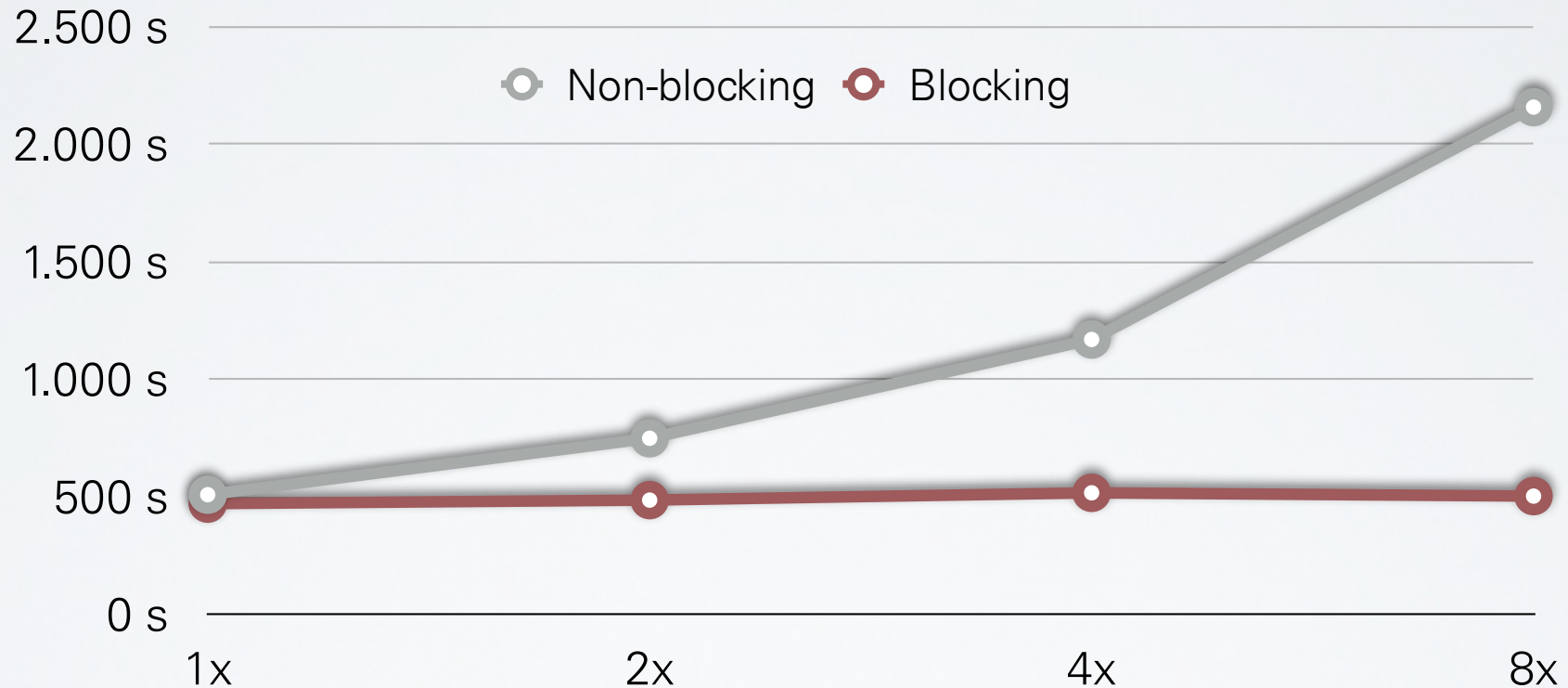
- Taurus 16 nodes w/ 16 Xeon E5-2690 (Sandy Bridge) @ 2.90GHz
- 1x - 8x oversubscription (256 - 2048 MPI ranks, same problem size)



Oversubscription factor (more ranks)

Application: COSMO-SPECS+FD4 (no load balancing)

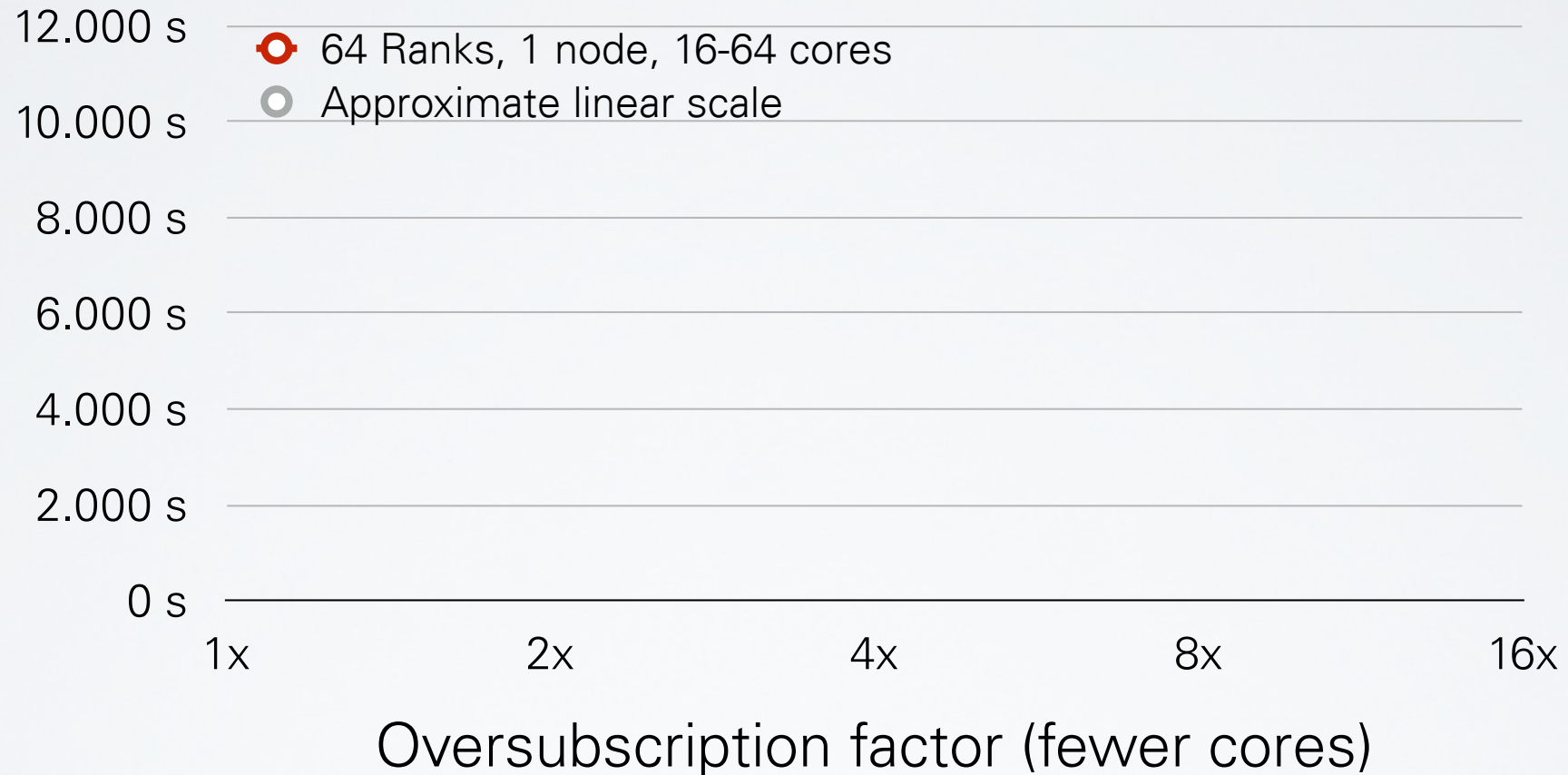
- Taurus 16 nodes w/ 16 Xeon E5-2690 (Sandy Bridge) @ 2.90GHz
- 1x - 8x oversubscription (256 - 2048 MPI ranks, same problem size)



Oversubscription factor (more ranks)

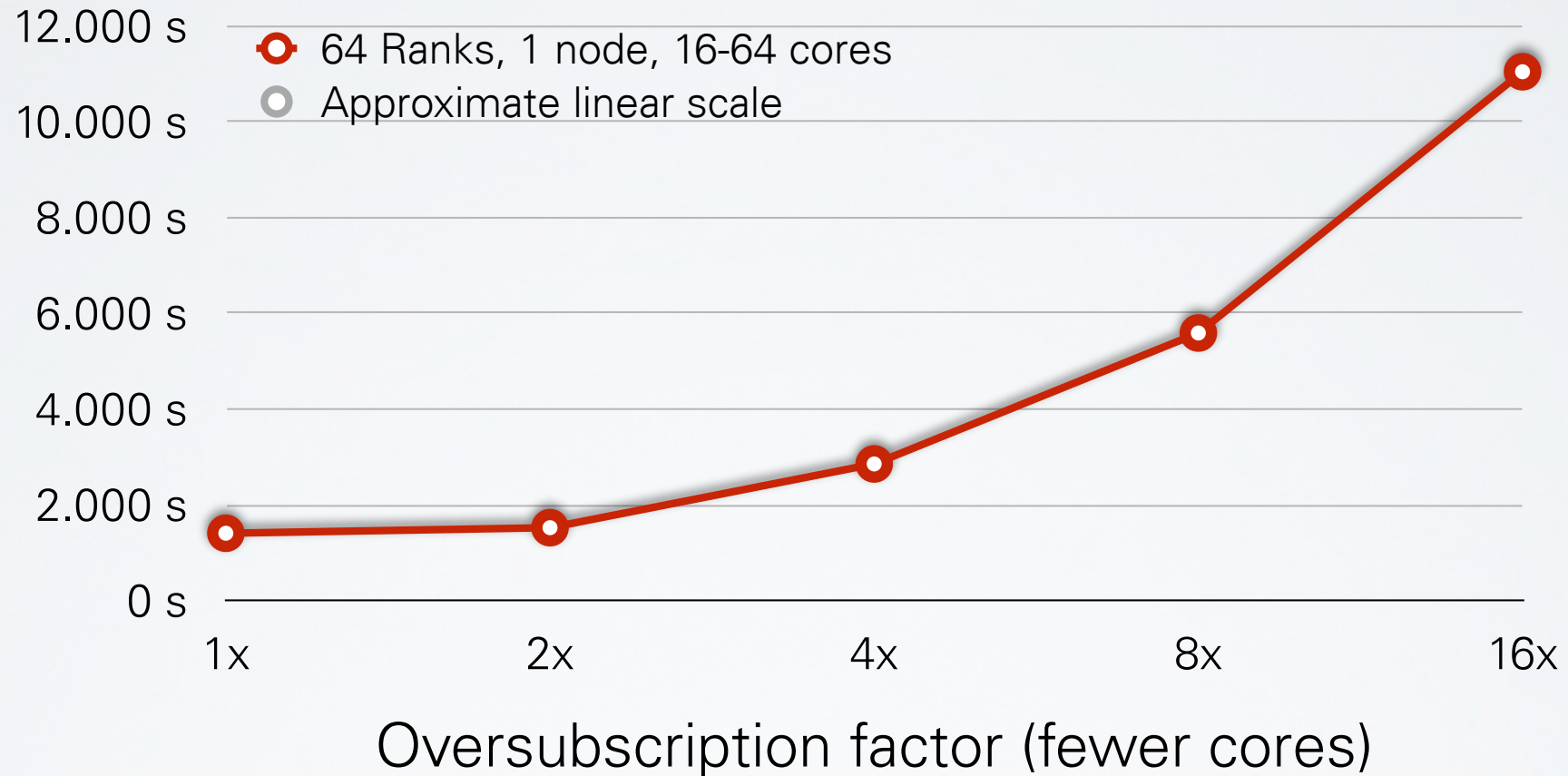
Application: COSMO-SPECS+FD4 (no load balancing)

- Taurus 16 nodes w/ 16 Xeon E5-2690 (Sandy Bridge) @ 2.90GHz
- 1x - 8x oversubscription (256 - 2048 MPI ranks, same problem size)



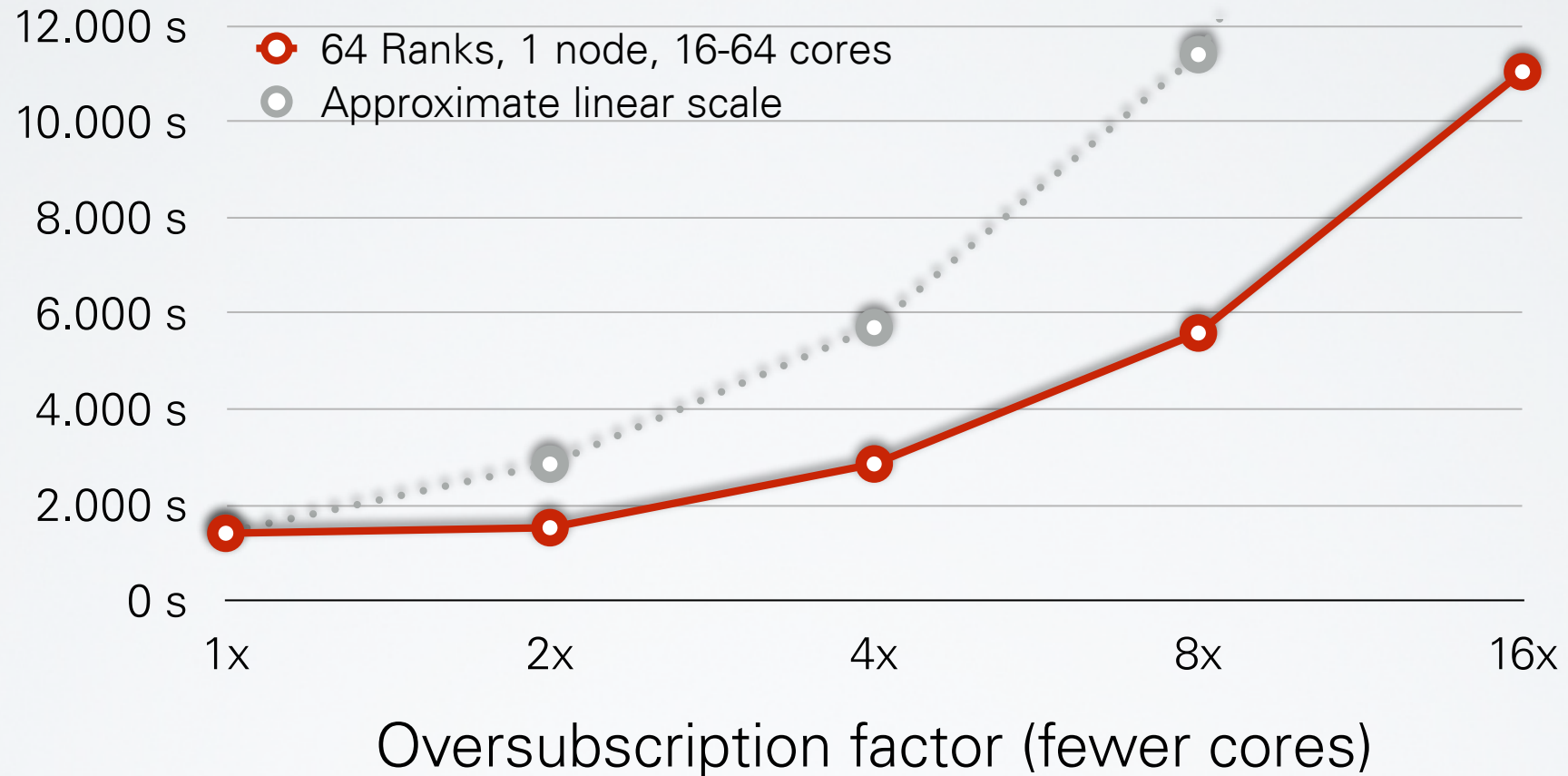
Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



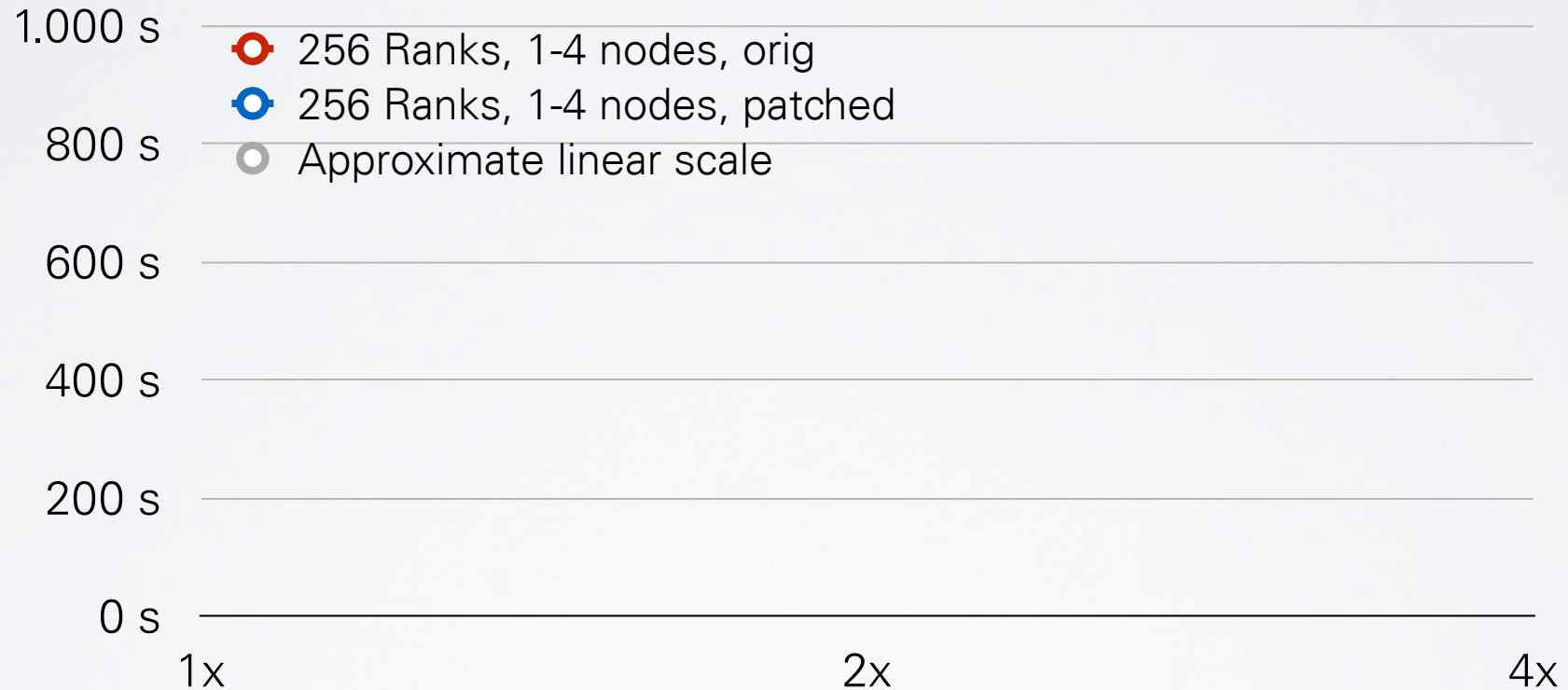
Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



Application: COSMO-SPECS+FD4 (no load balancing)

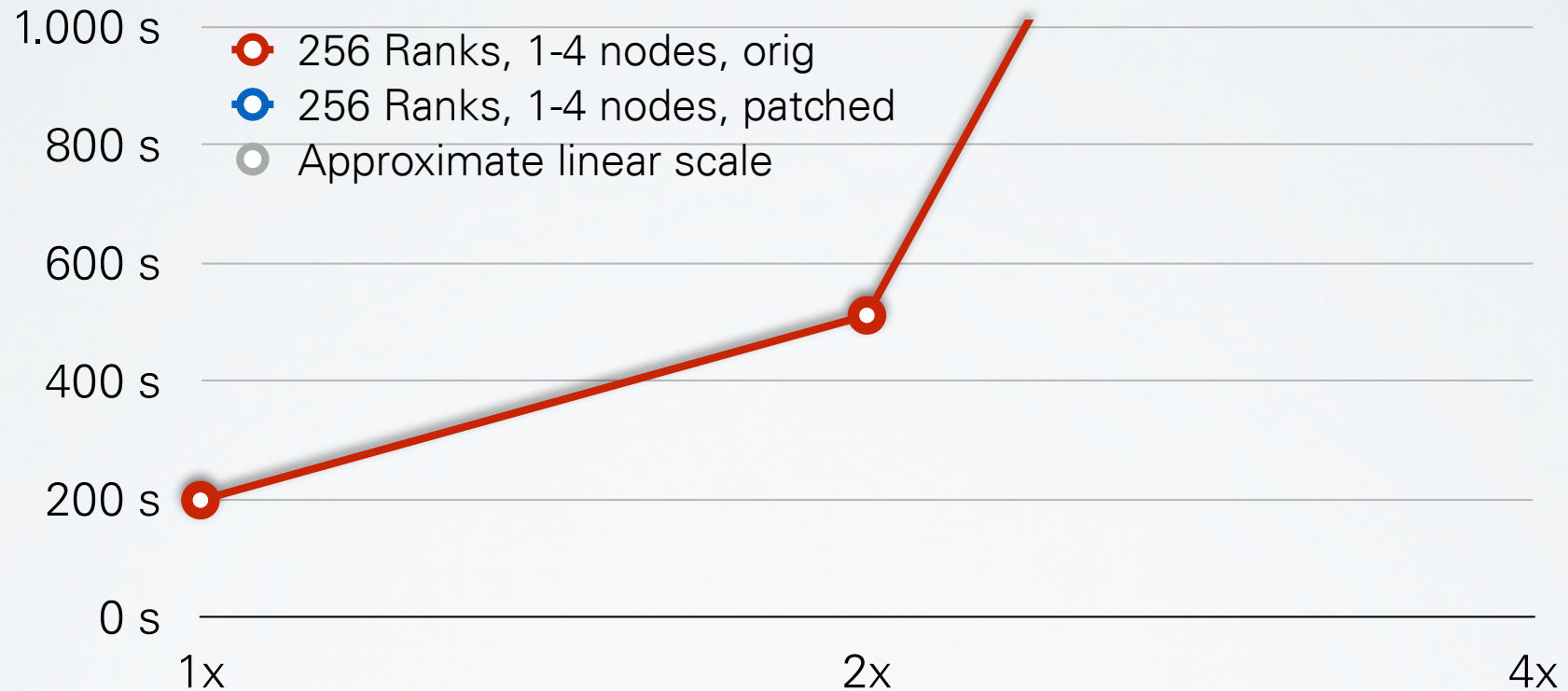
- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



Oversubscription factor (fewer cores)

Application: COSMO-SPECS+FD4 (no load balancing)

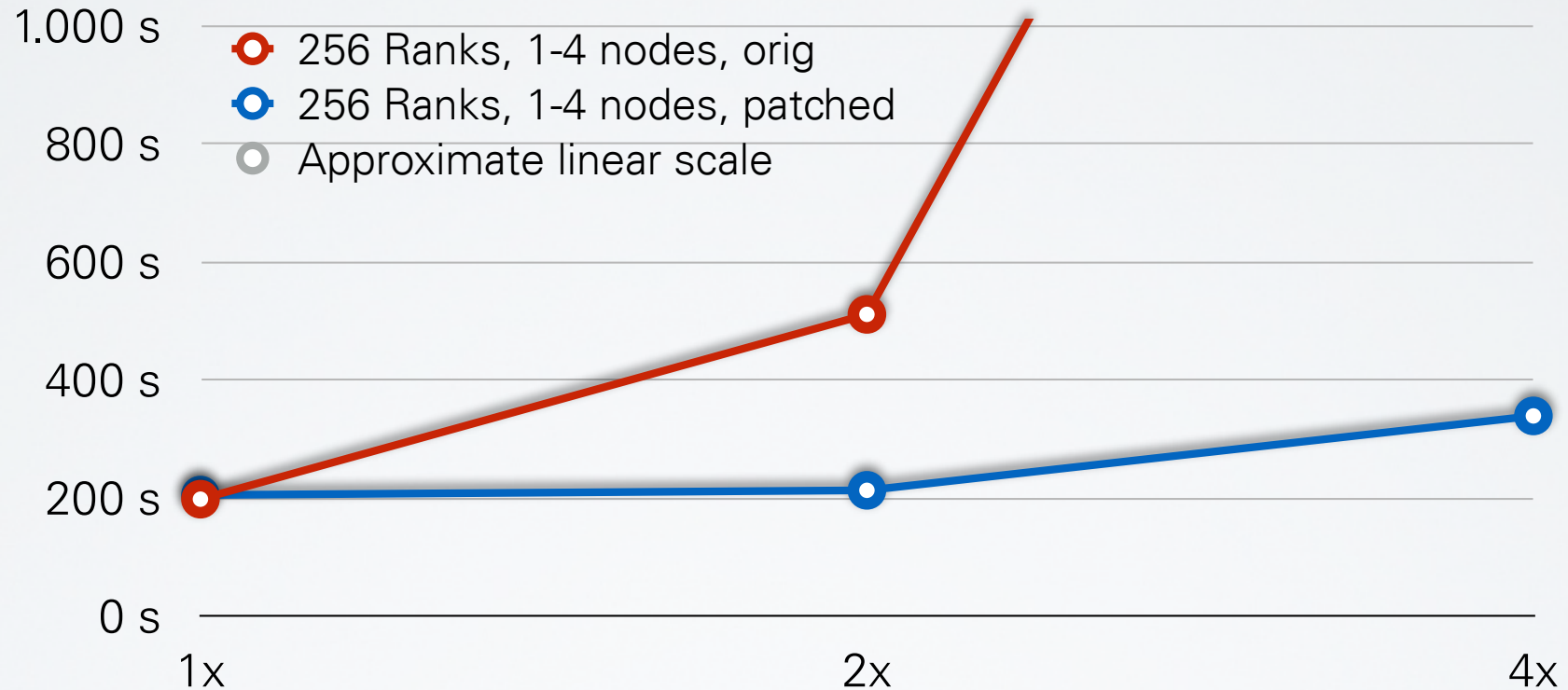
- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



Oversubscription factor (fewer cores)

Application: COSMO-SPECS+FD4 (no load balancing)

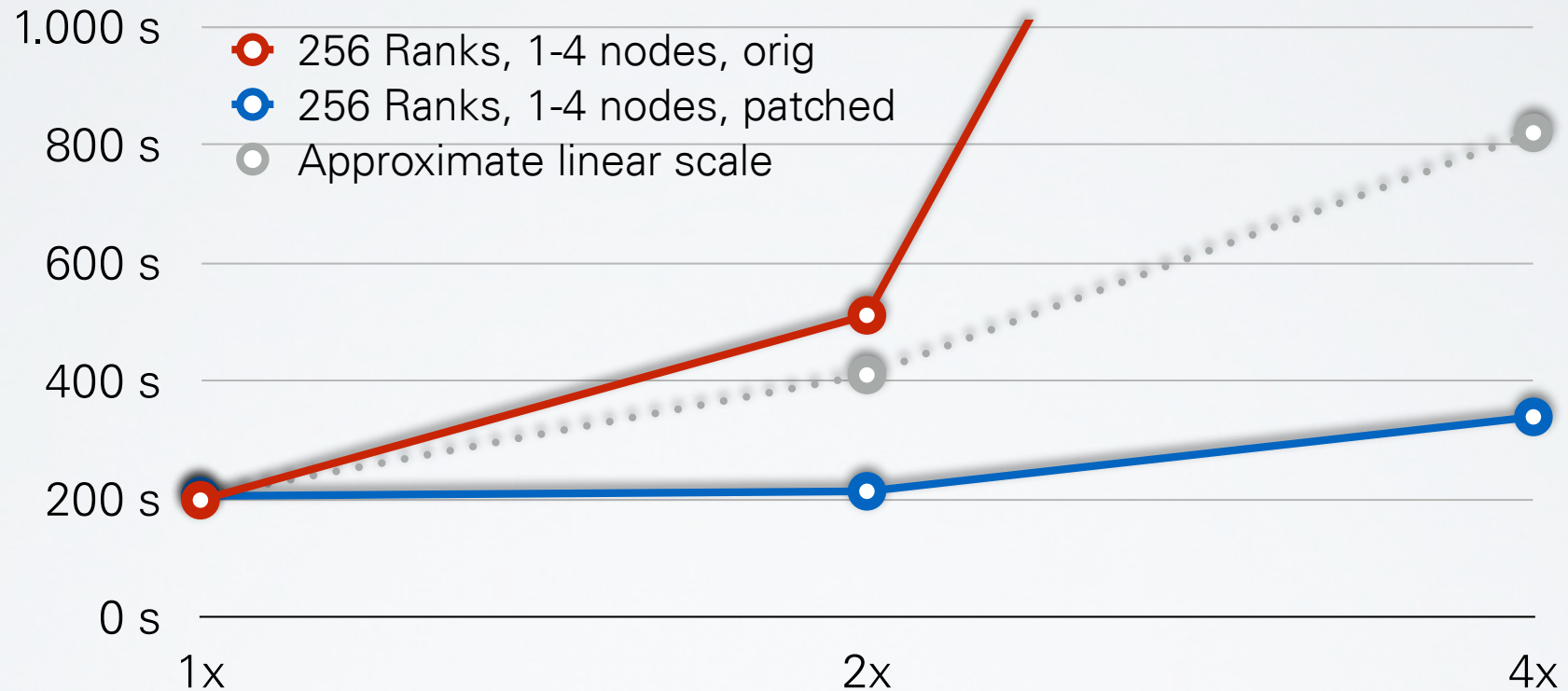
- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



Oversubscription factor (fewer cores)

Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



Oversubscription factor (fewer cores)

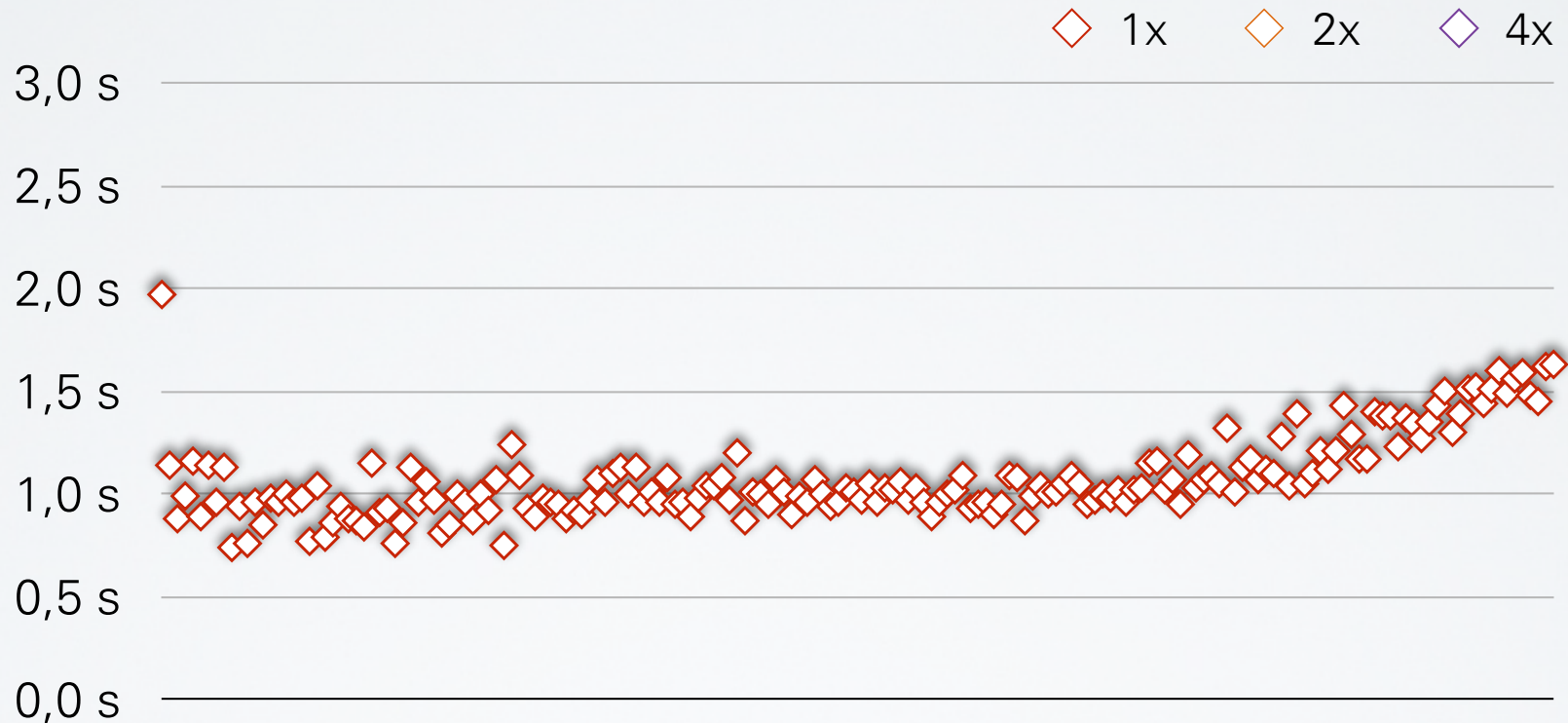
Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



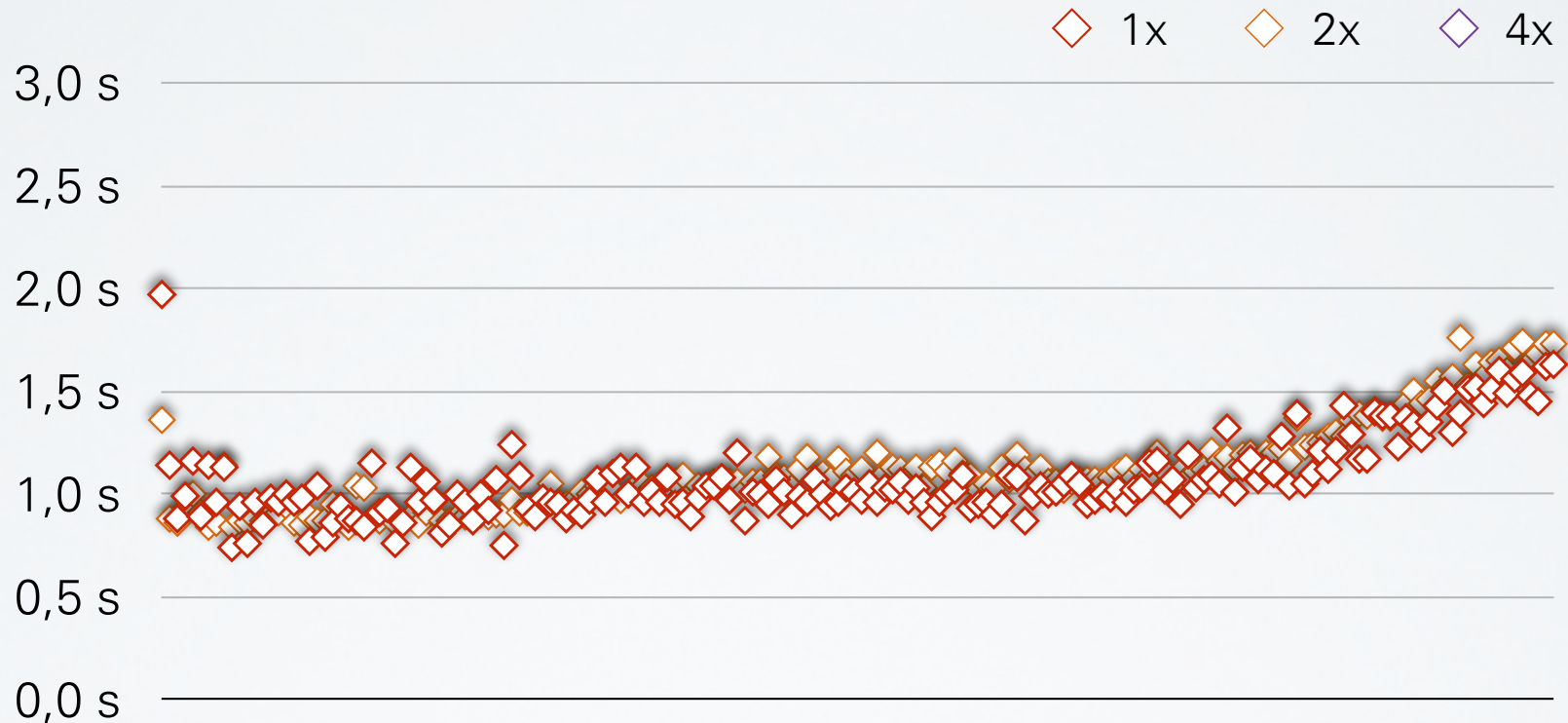
Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



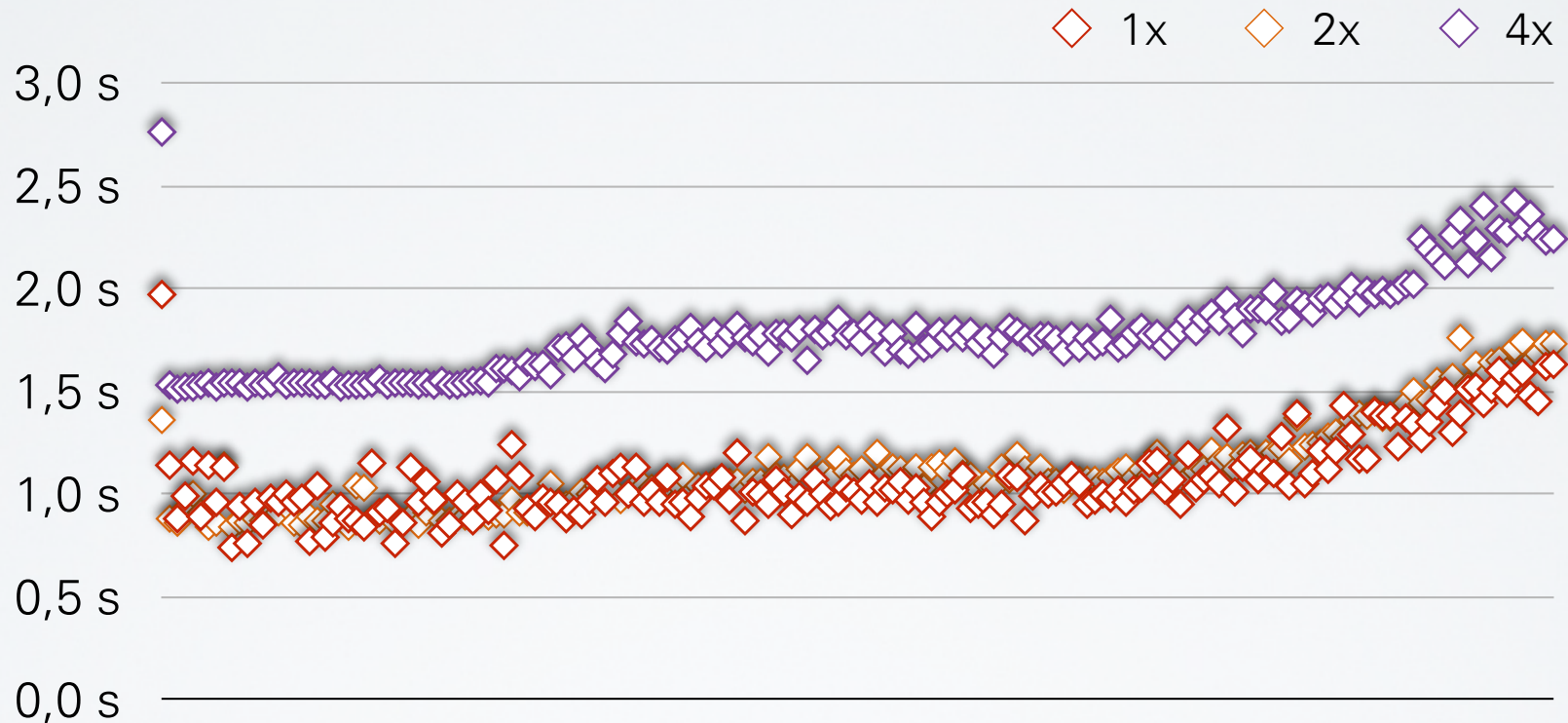
Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



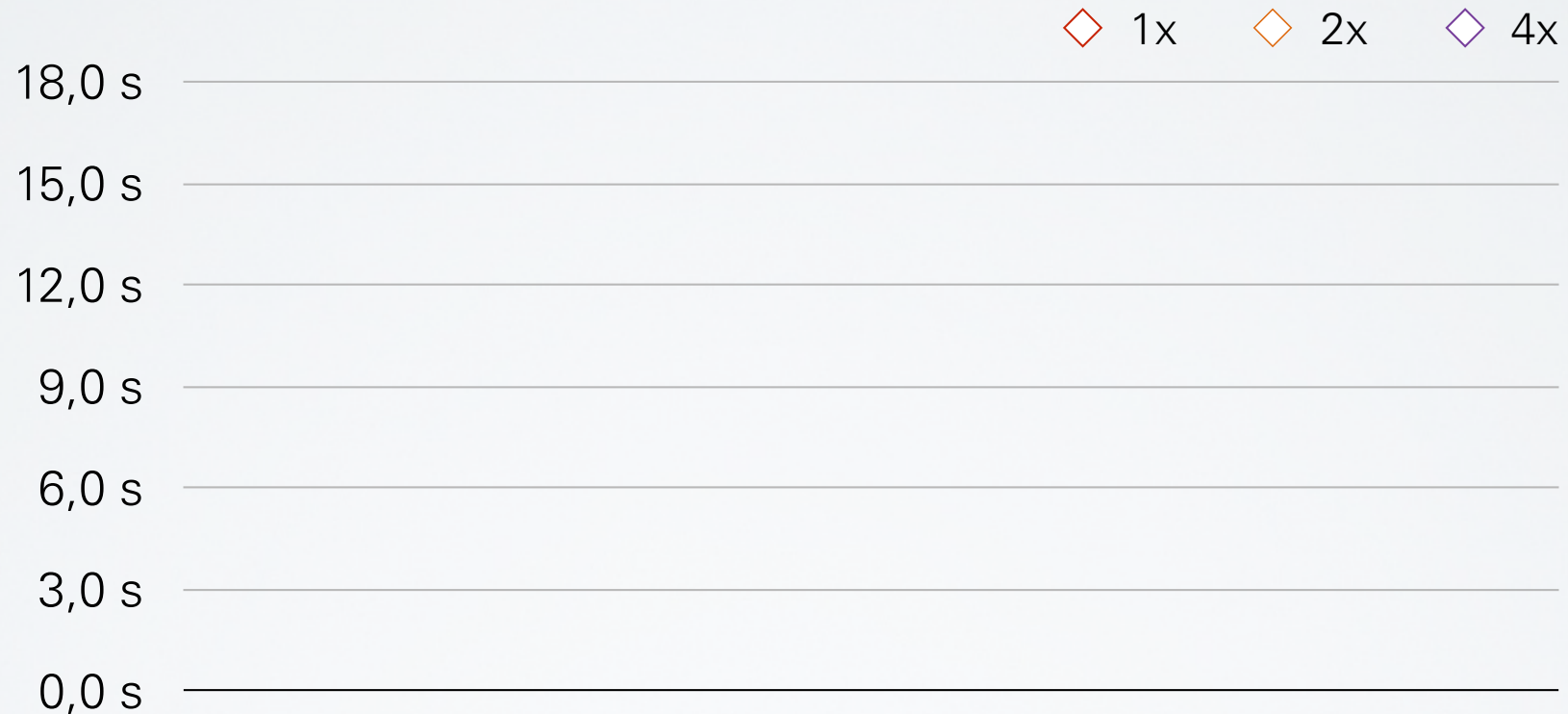
Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



Application: COSMO-SPECS+FD4 (no load balancing)

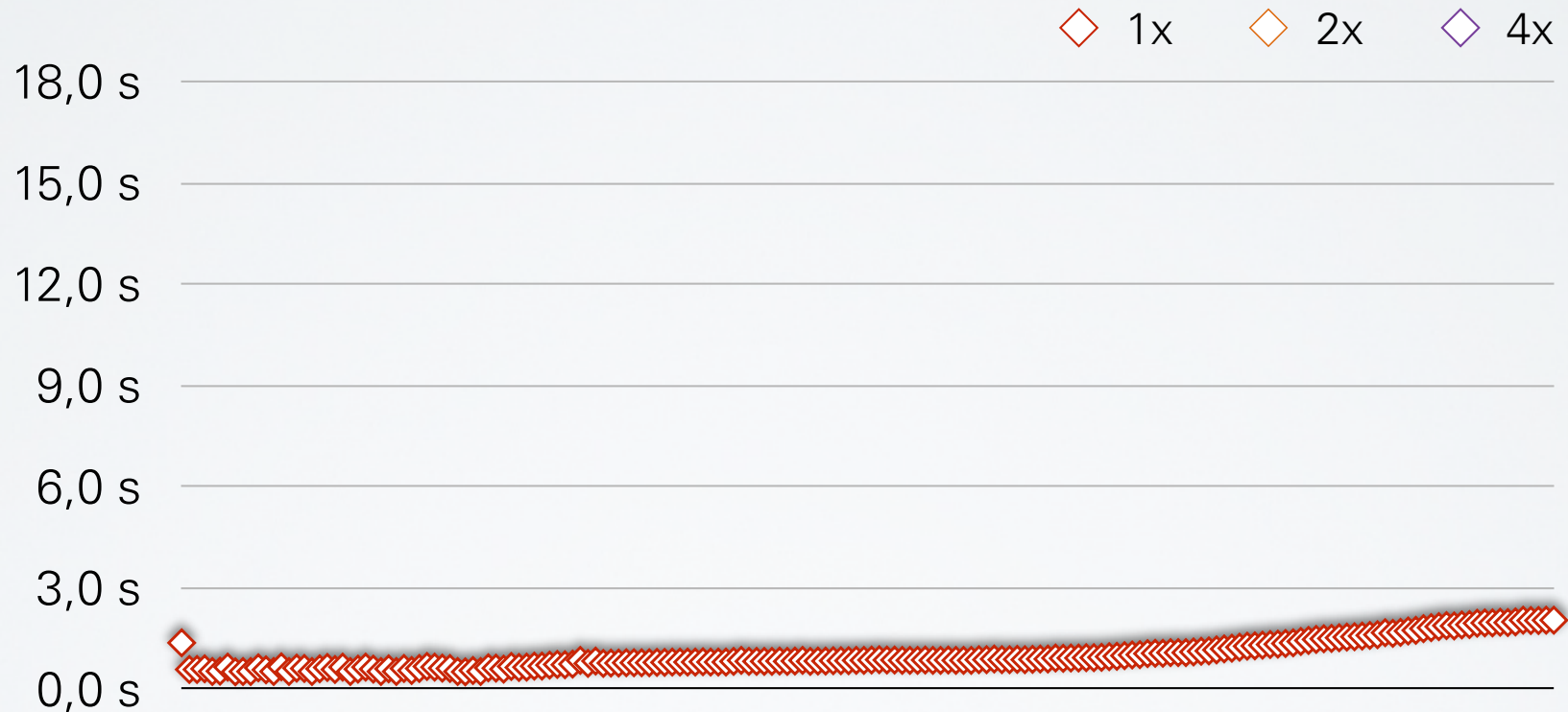
- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced



Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced

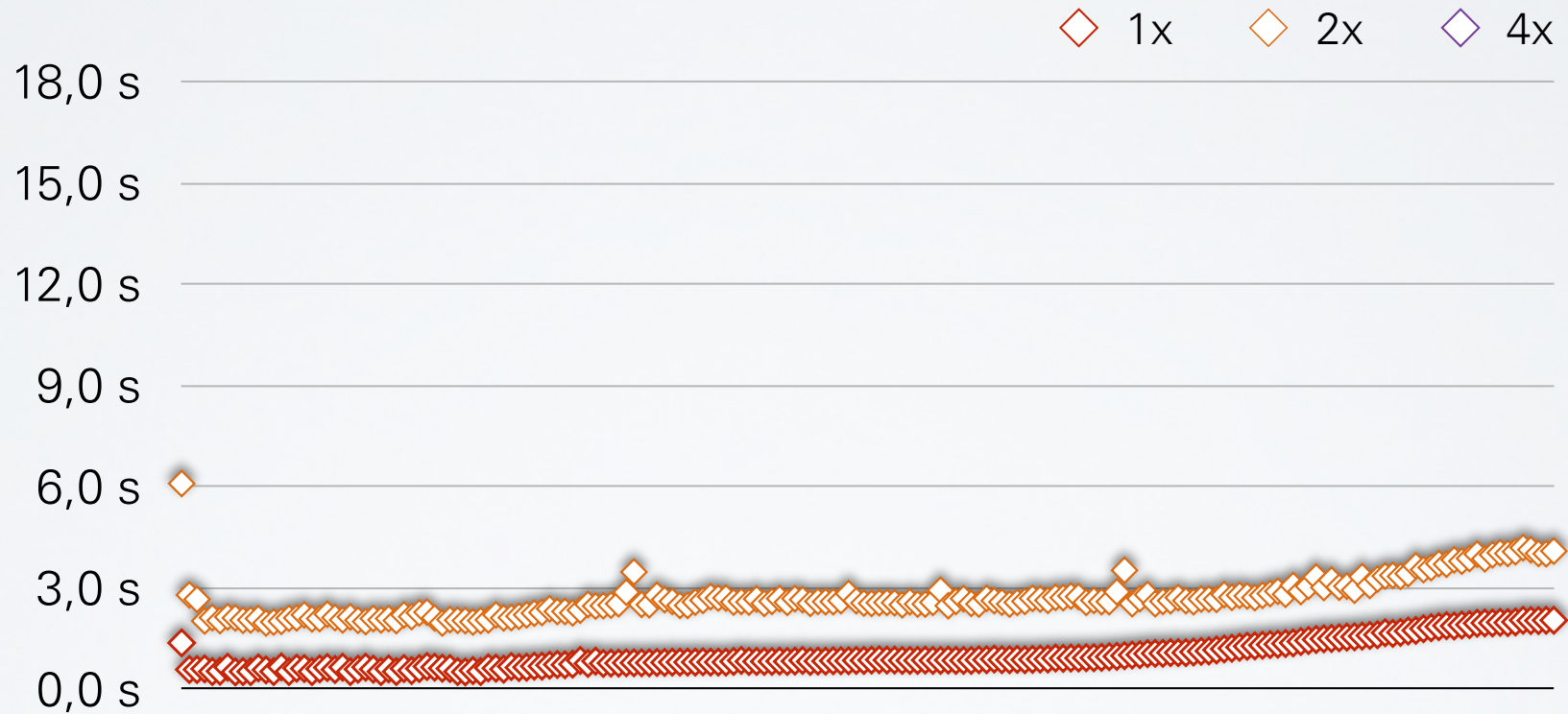
ORIG: STEP TIME



Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced

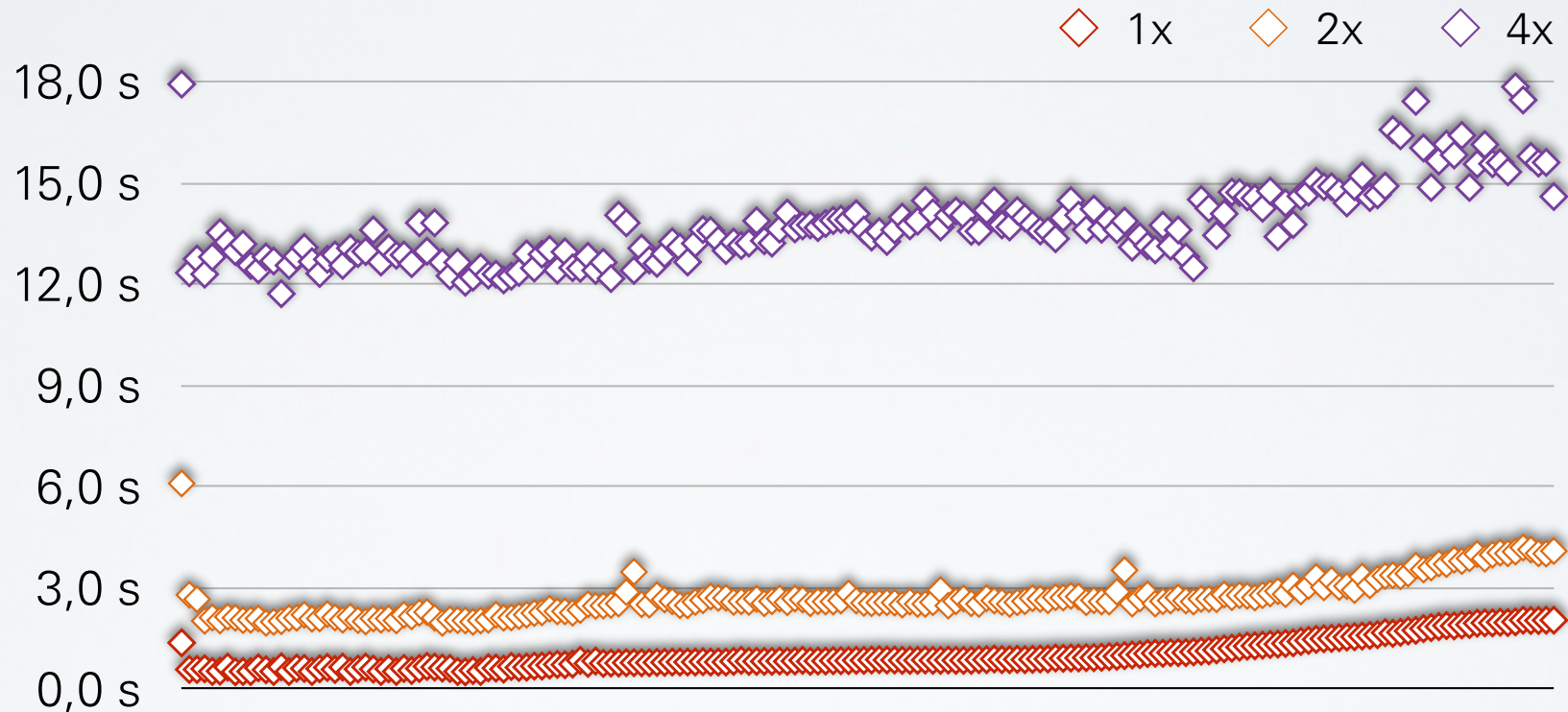
ORIG: STEP TIME



Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced

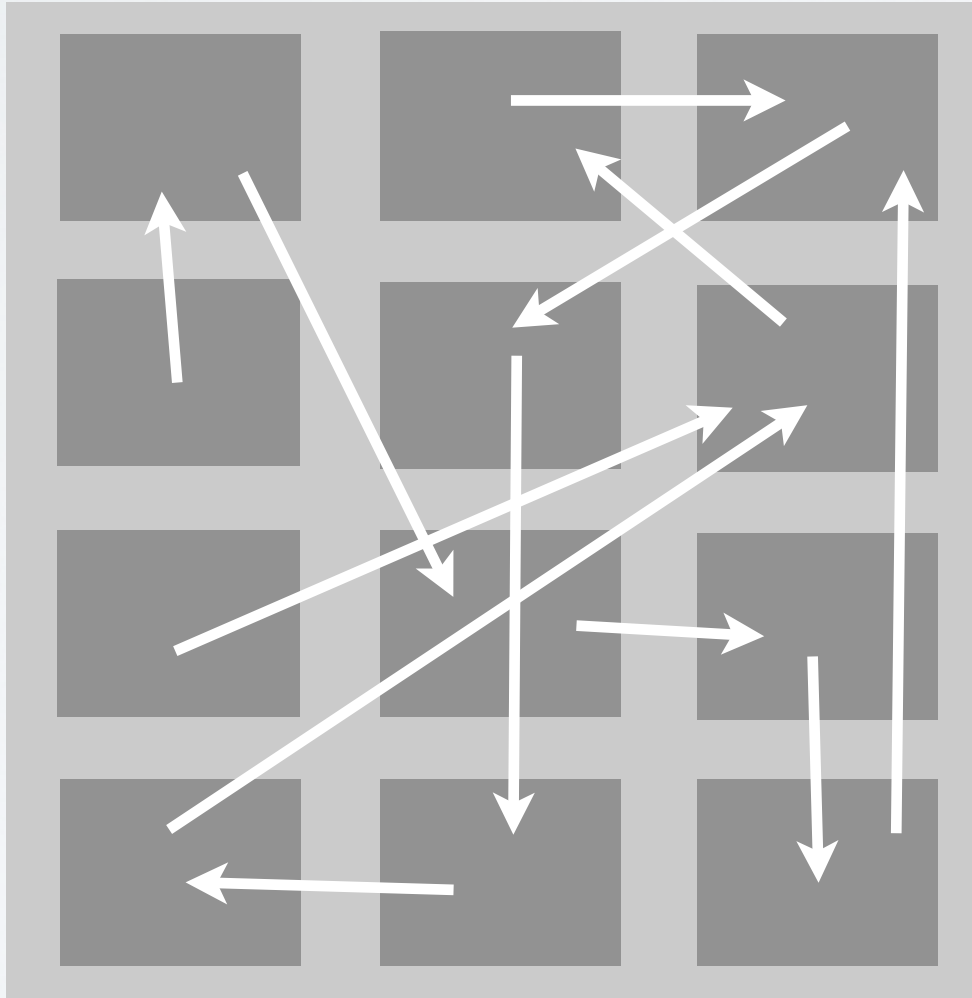
ORIG: STEP TIME



Application: COSMO-SPECS+FD4 (no load balancing)

- ATLAS nodes w/ 64 AMD Opteron 6274 cores @ 2.2 GHz
- Number of ranks remained constant, but number of cores was reduced

EXPERIMENTS: GOSSIP SCALABILITY

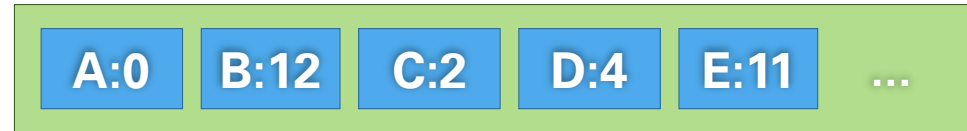


Distributed Bulletin Board

- Each node keeps vector with per-node info (own + info received from others)
- Once per time step, each node sends to 1 other randomly selected node a subset of its own vector entries (called "window")
- Node merges received window entries into local vector (if newer)

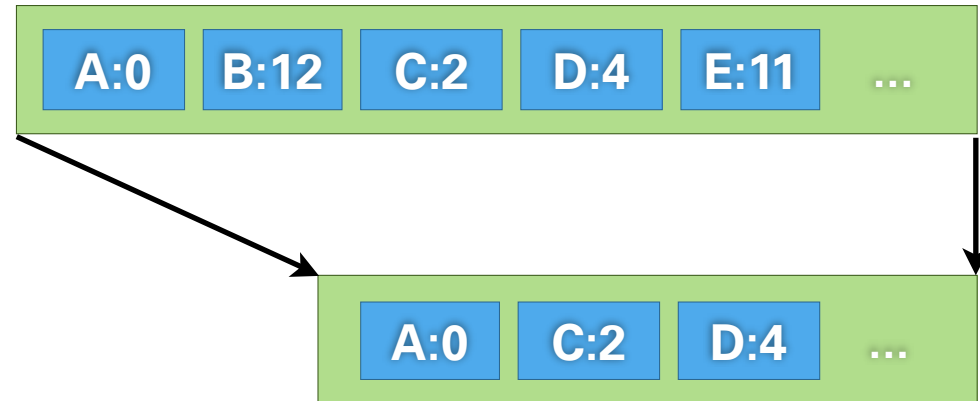
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



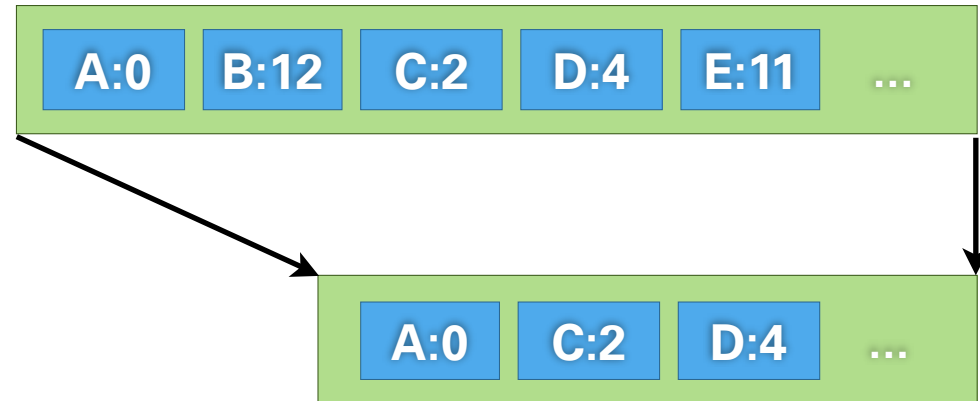
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



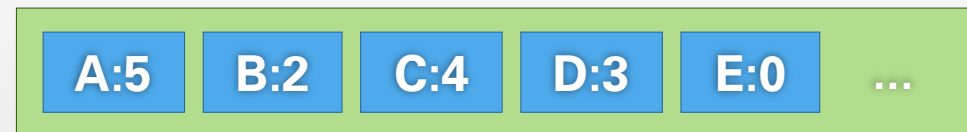
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



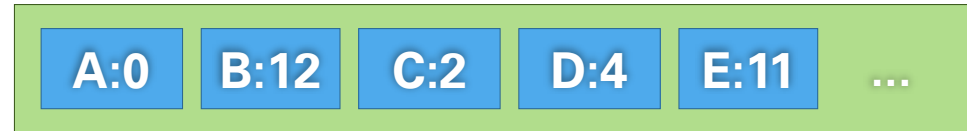
Upon receiving a window:

- Update the received entries' age (+1 for transfer)
- Update entries in local vector where newer information has been received



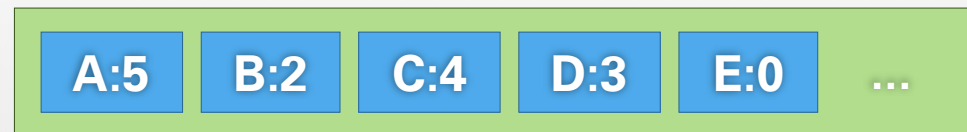
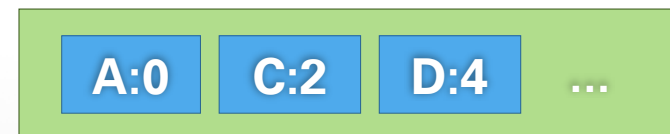
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



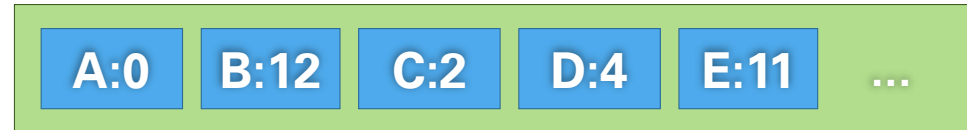
Upon receiving a window:

- Update the received entries' age (+1 for transfer)
- Update entries in local vector where newer information has been received



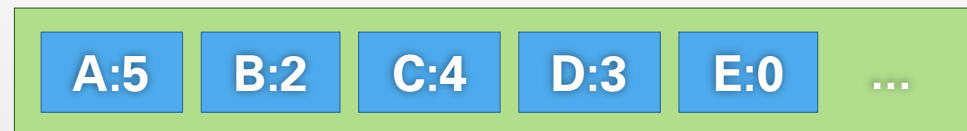
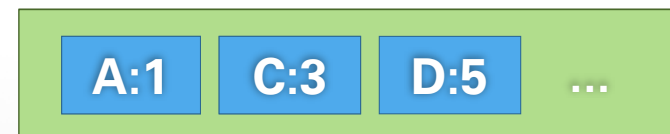
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



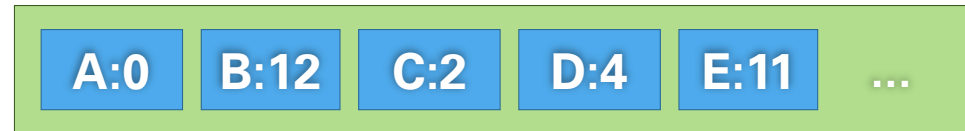
Upon receiving a window:

- Update the received entries' age (+1 for transfer)
- Update entries in local vector where newer information has been received



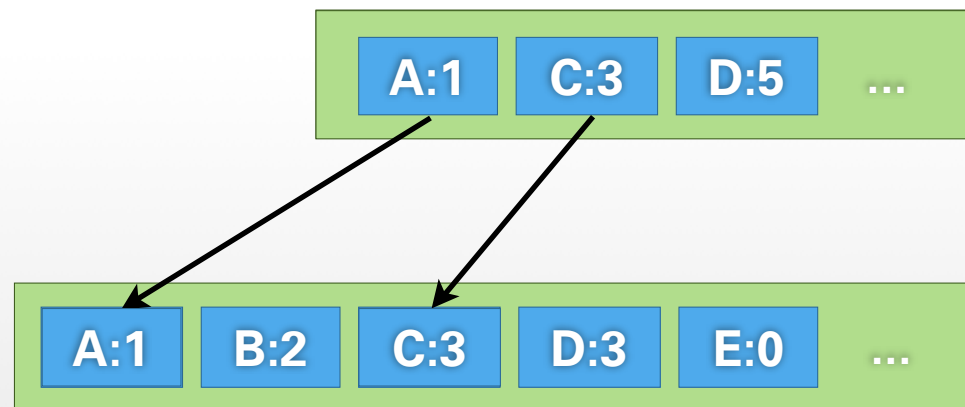
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node

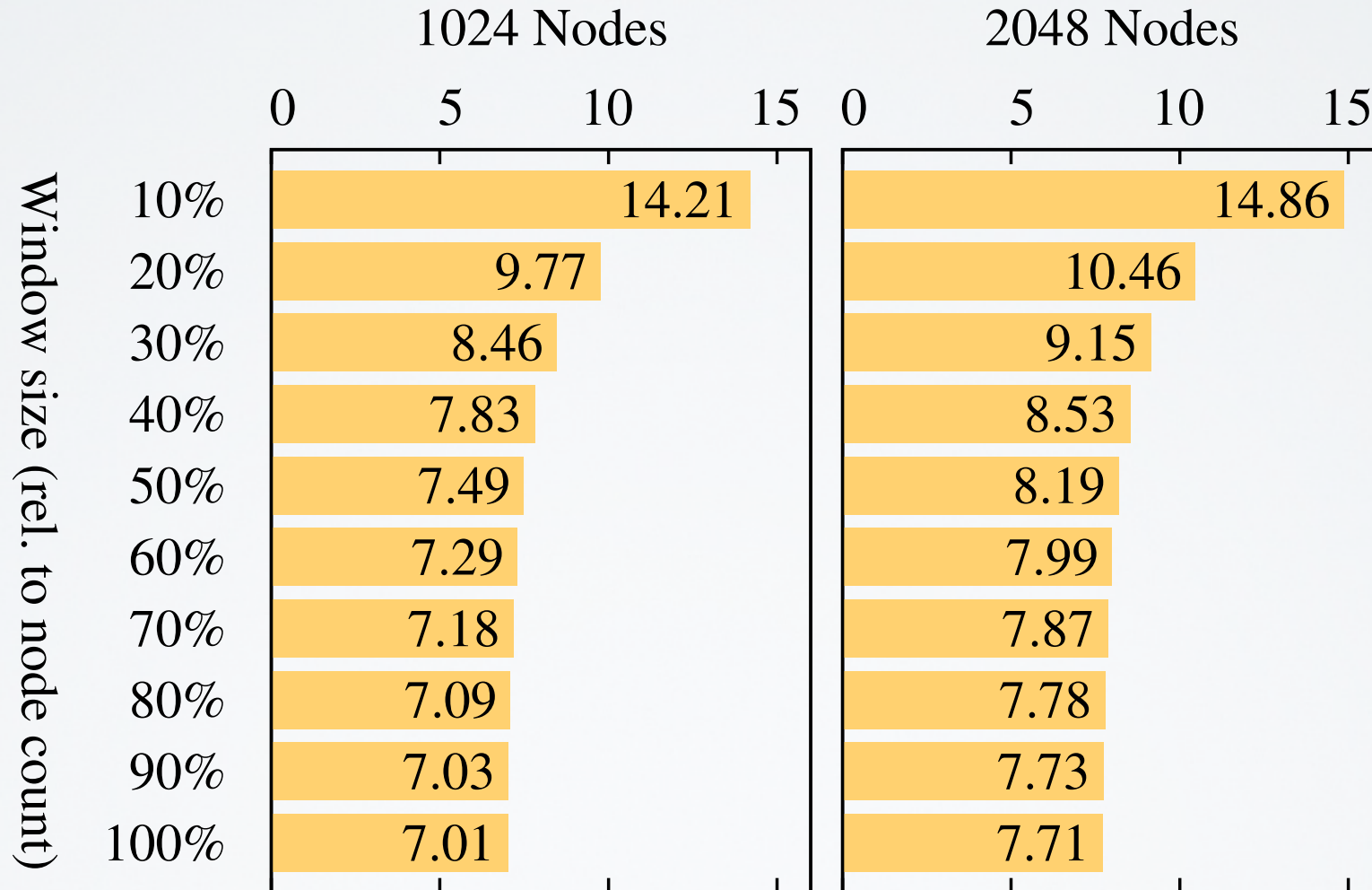


Upon receiving a window:

- Update the received entries' age (+1 for transfer)
- Update entries in local vector where newer information has been received

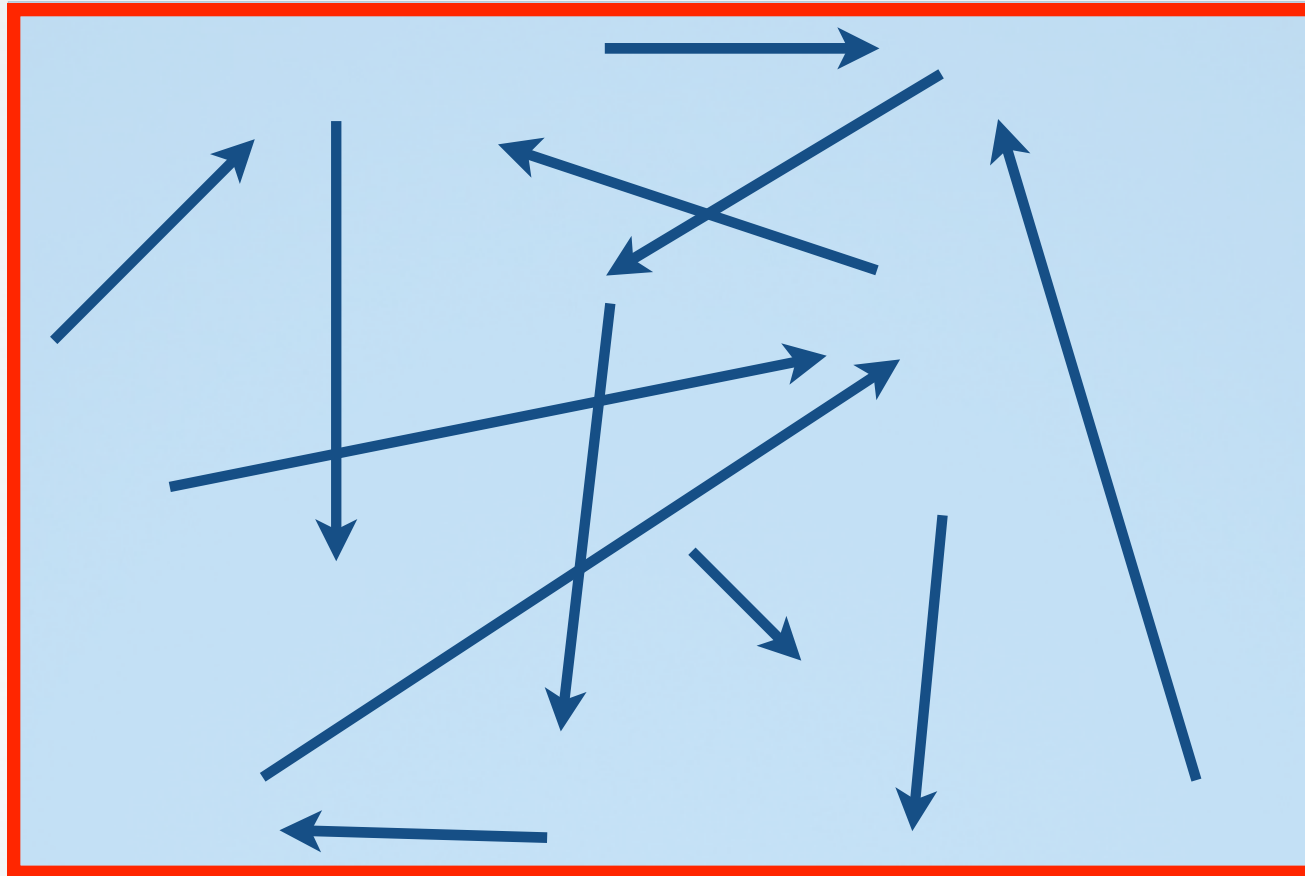


WINDOW SIZE

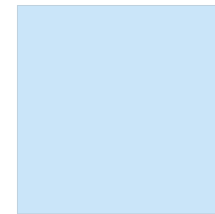


Colony nodes	Method	Circulating among colony nodes					
		windows not exceeding age					whole vector
		2	4	6	8	10	
128	Approx.	19.15	6.00	4.93	4.89	4.89	4.89
	Simulation	18.87	6.04	4.97	4.92	4.95	
	Measured	18.75	5.99	4.94	4.88	4.90	
256	Approx.	36.49	8.49	5.70	5.57	5.57	5.57
	Simulation	36.33	8.57	5.77	5.63	5.62	
	Measured	36.06	8.55	5.77	5.60	5.60	
512	Approx.	71.15	13.27	6.70	6.26	6.25	6.25
	Simulation	71.01	13.34	6.81	6.34	6.32	
	Measured	70.85	13.37	6.78	6.31	6.28	
1K	Approx.	140.44	22.69	8.21	6.99	6.94	6.94
	Simulation	139.76	22.73	8.33	7.06	7.01	
	Measured	140.14	22.83	8.32	7.04	6.98	
2K	Approx.	279.03	41.47	10.90	7.79	7.63	7.63
	Simulation	267.82	41.58	11.08	7.89	7.71	
	Measured	278.94	41.66	11.03	7.84	7.66	
4K	Approx.	556.20	78.99	16.06	8.83	8.34	8.32
	Simulation	479.96	79.10	16.23	8.95	8.42	
	Measured	556.20	79.39	16.24	8.87	8.33	
8K	Approx.	1,110.53	154.02	26.26	10.44	9.07	9.01
	Simulation	798.97	153.80	26.48	10.59	9.43	
	Measured	1,102.99	155.16	26.51	10.44	8.98	
1M	Approx.	141,911	19,209	2,605	360	58	13.86
1G	Approx.	145M	19M	2M	360K	48K	20.79

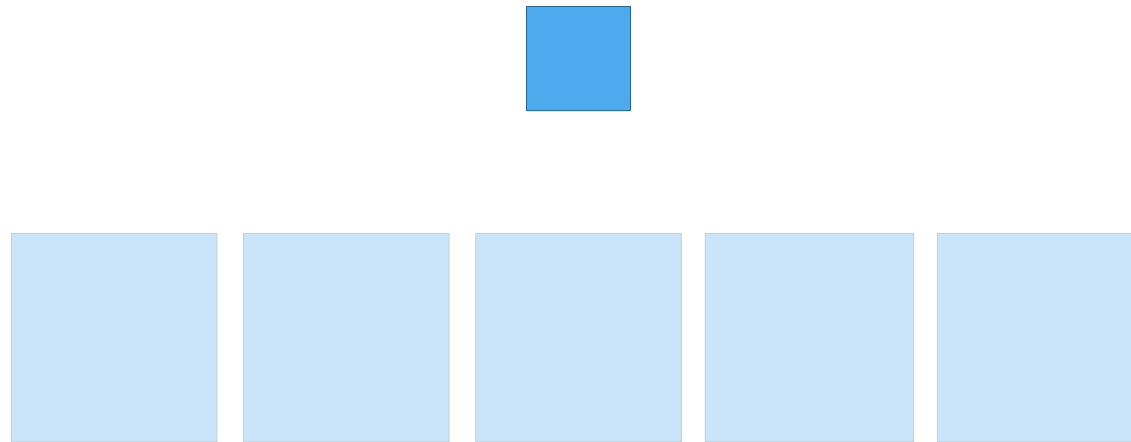
Colony nodes	Method	Circulating among colony nodes					
		windows not exceeding age					whole vector
		2	4	6	8	10	
128	Approx.	19.15	6.00	4.93	4.89	4.89	4.89
	Simulation	18.87	6.04	4.97	4.92	4.95	
	Measured	18.75	5.99	4.94	4.88	4.90	
256	Approx.	36.49	8.49	5.70	5.57	5.57	5.57
	Simulation	36.33	8.57	5.77	5.63	5.62	
	Measured	36.06	8.55	5.77	5.60	5.60	
512	Approx.	71.15	13.27	6.70	6.26	6.25	6.25
	Simulation	71.01	13.34	6.81	6.34	6.32	
	Measured	70.85	13.37	6.78	6.31	6.28	
1K	Approx.	140.44	22.69	8.21	6.99	6.94	6.94
	Simulation	139.76	22.73	8.33	7.06	7.01	
	Measured	140.14	22.83	8.32	7.04	6.98	
2K	Approx.	279.03	41.47	10.90	7.79	7.63	7.63
	Simulation	267.82	41.58	11.08	7.89	7.71	
	Measured	278.94	41.66	11.03	7.84	7.66	
4K	Approx.	556.20	78.99	16.06	8.83	8.34	8.32
	Simulation	479.96	79.10	16.23	8.95	8.42	
	Measured	556.20	79.39	16.24	8.87	8.33	
8K	Approx.	1,110.53	154.02	26.26	10.44	9.07	9.01
	Simulation	798.97	153.80	26.48	10.59	9.43	
	Measured	1,102.99	155.10	26.51	10.44	8.98	
1M	Approx.	141,911	19,209	2,605	360	58	13.86
1G	Approx.	145M	19M	2M	360K	48K	20.79

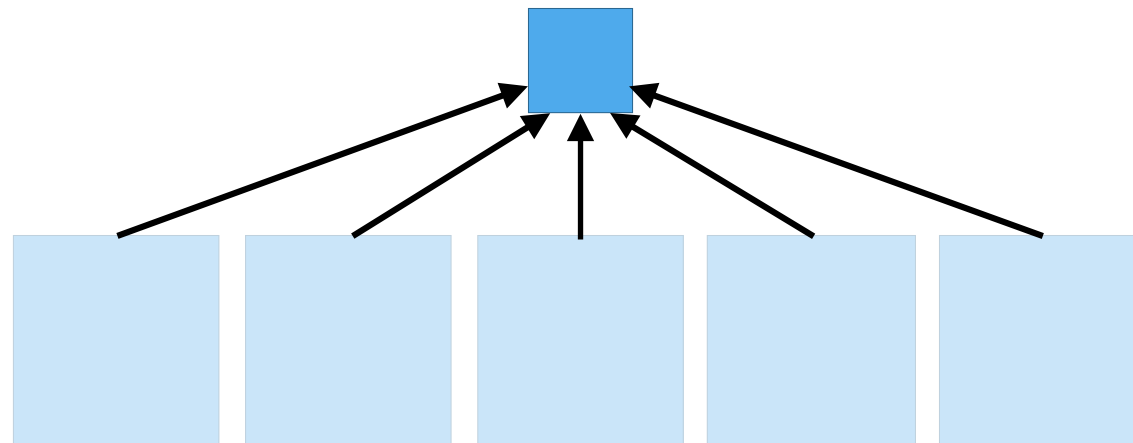


Problem: average age or window sizes too big for extreme numbers of nodes

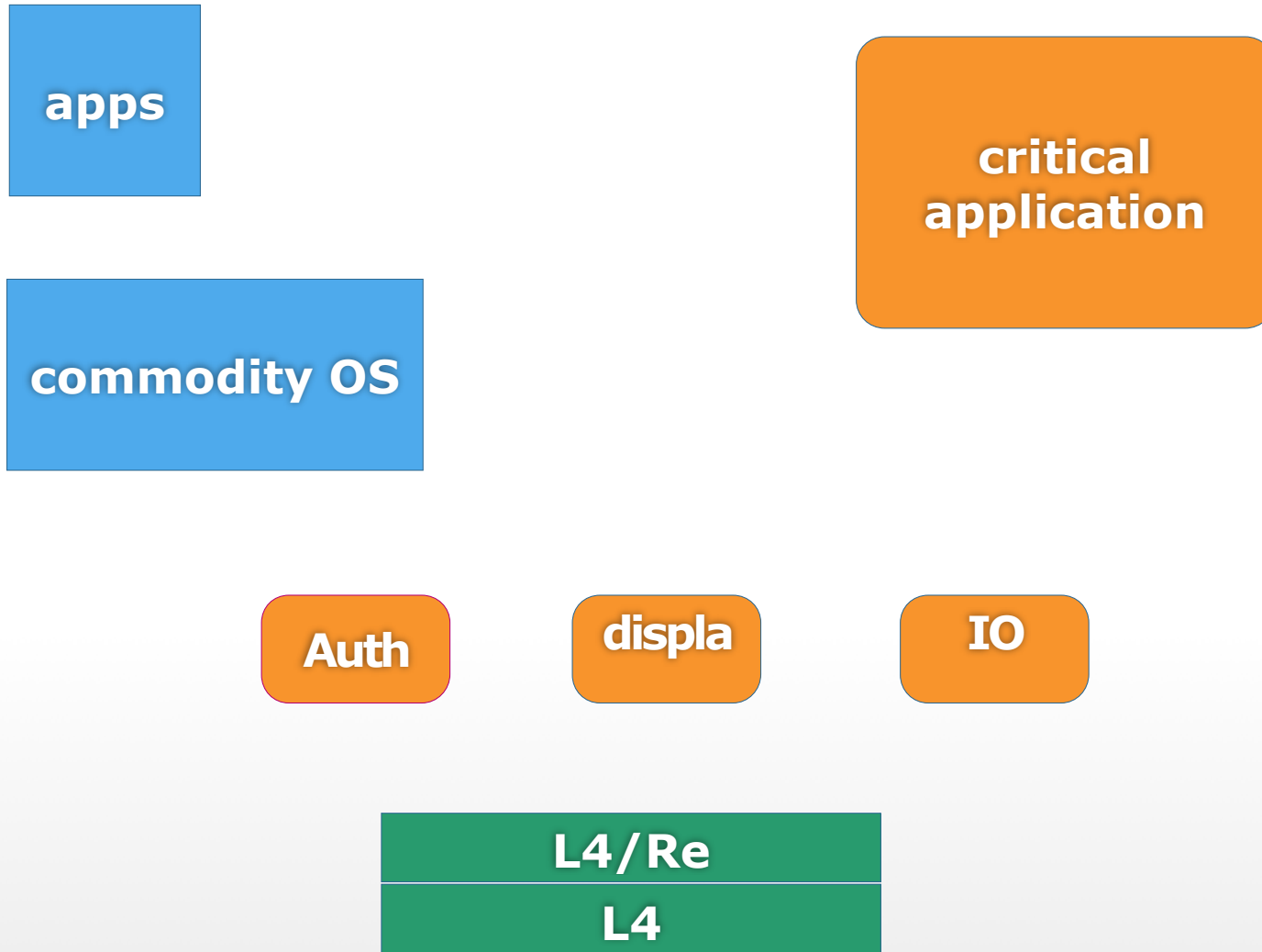








SYSTEM ARCHITECTURE

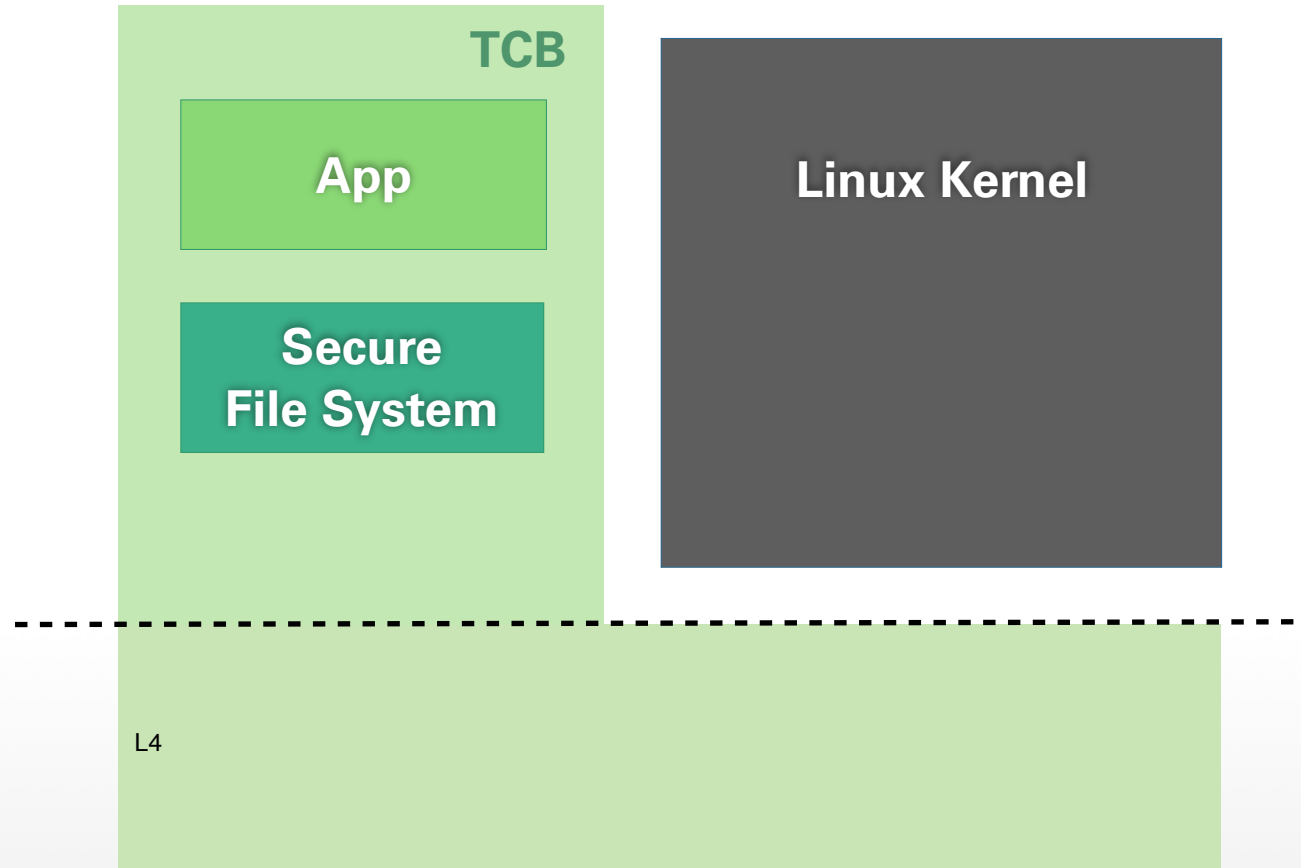


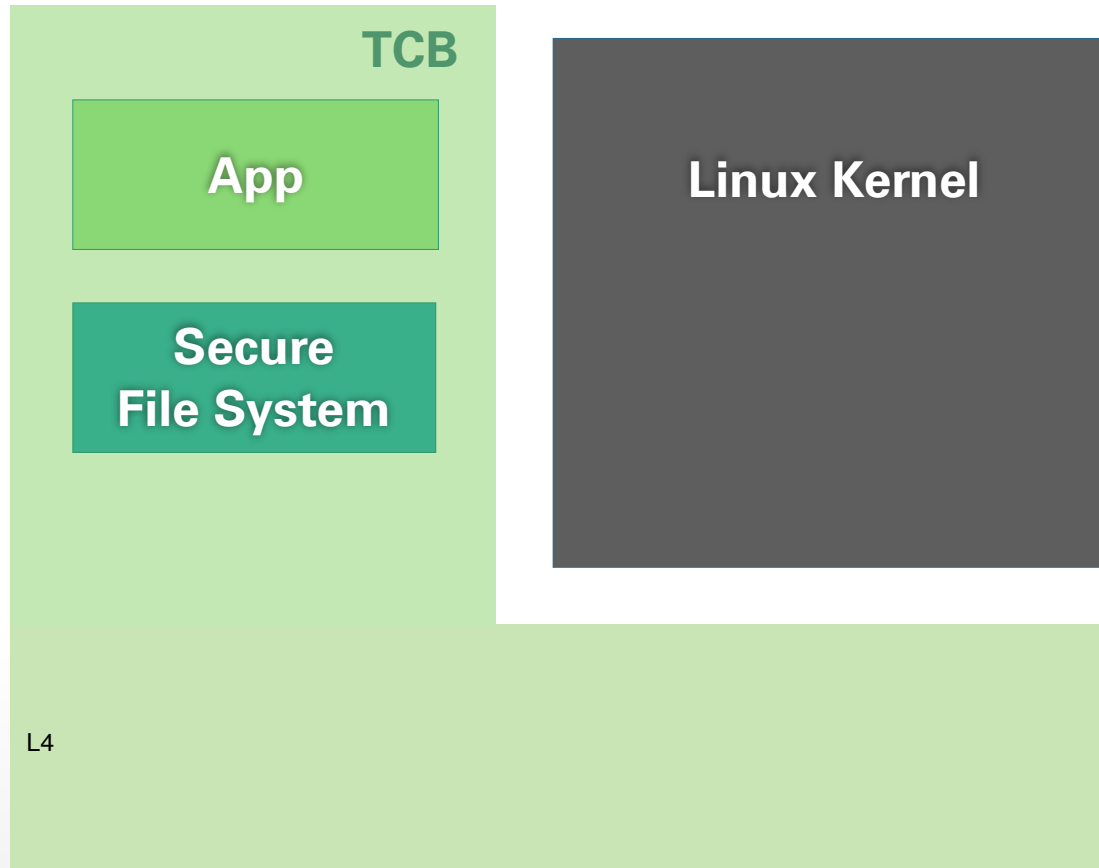


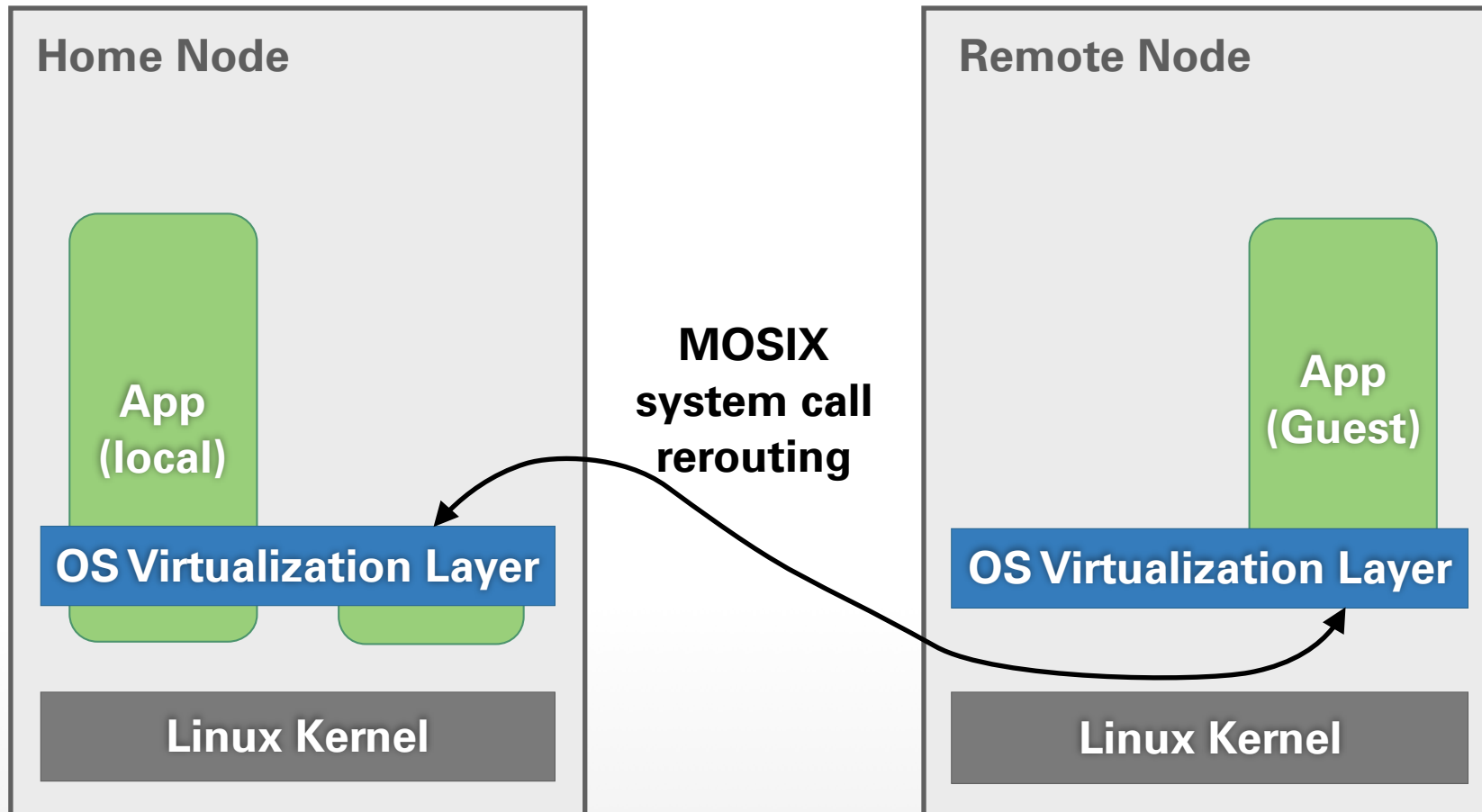


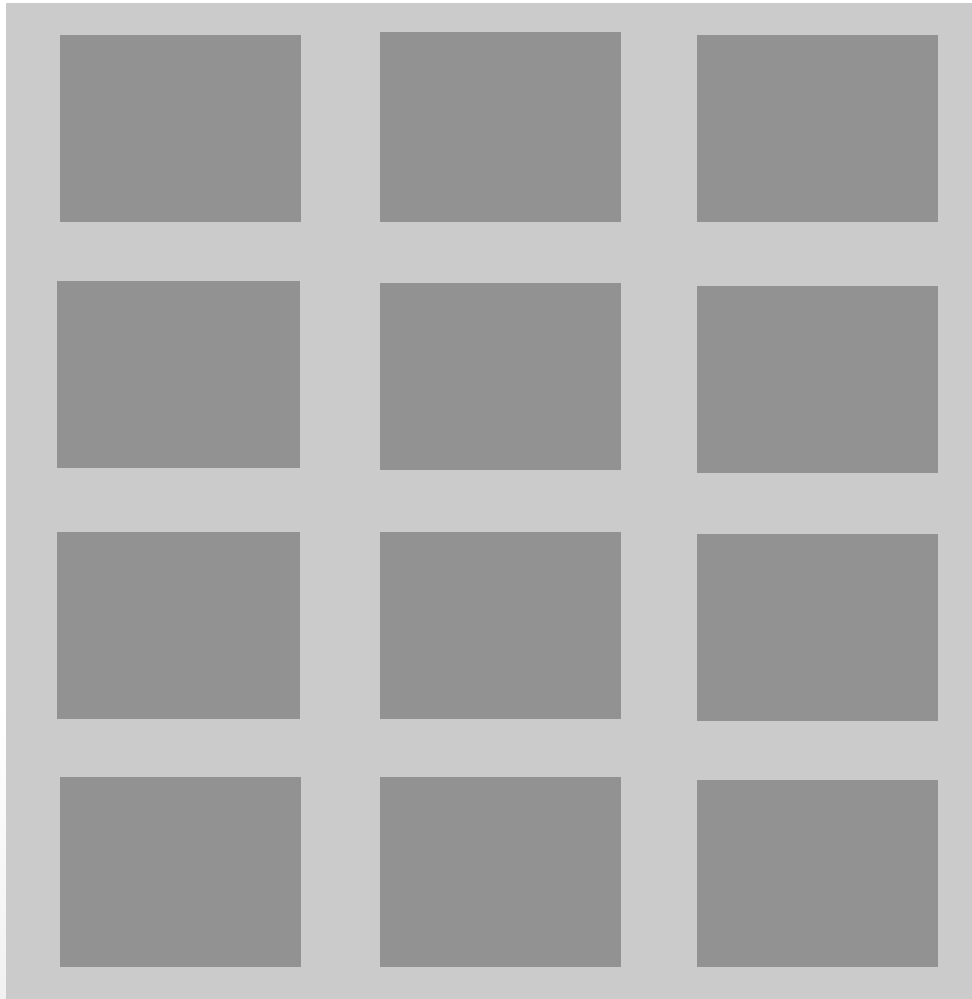
“ Merkel
Phone “





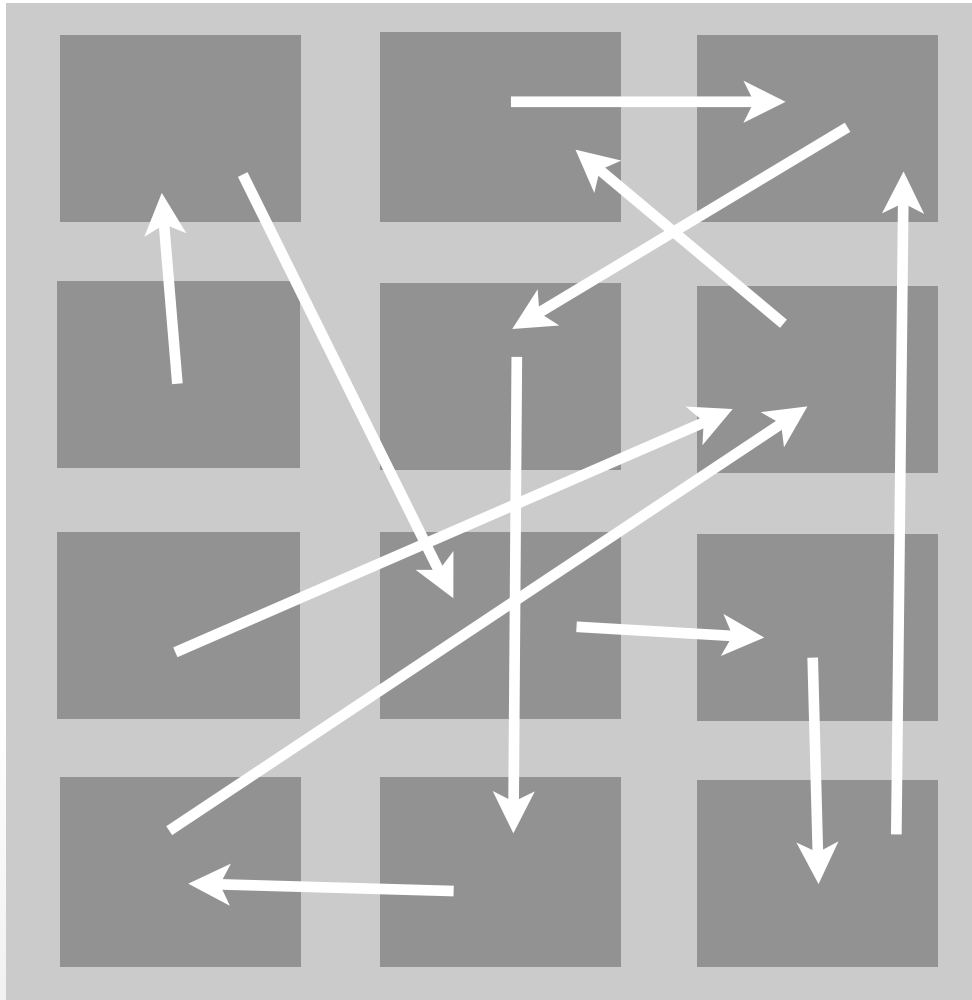






Distributed Bulletin Board

- Each node keeps vector with per-node info (own + info received from others)
- Once per time step, each node sends to 1 other randomly selected node a subset of its own vector entries (called “window”)
- Node merges received window entries into local vector (if newer)

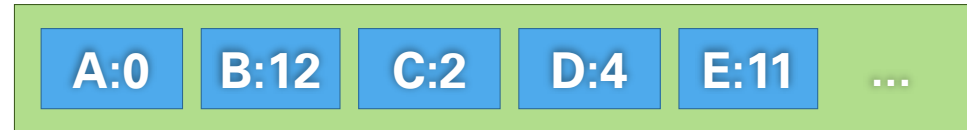


Distributed Bulletin Board

- Each node keeps vector with per-node info (own + info received from others)
- Once per time step, each node sends to 1 other randomly selected node a subset of its own vector entries (called "window")
- Node merges received window entries into local vector (if newer)

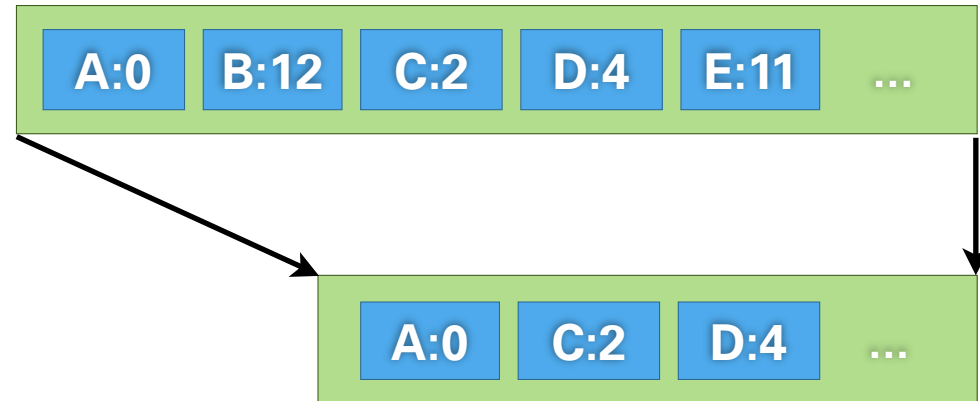
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



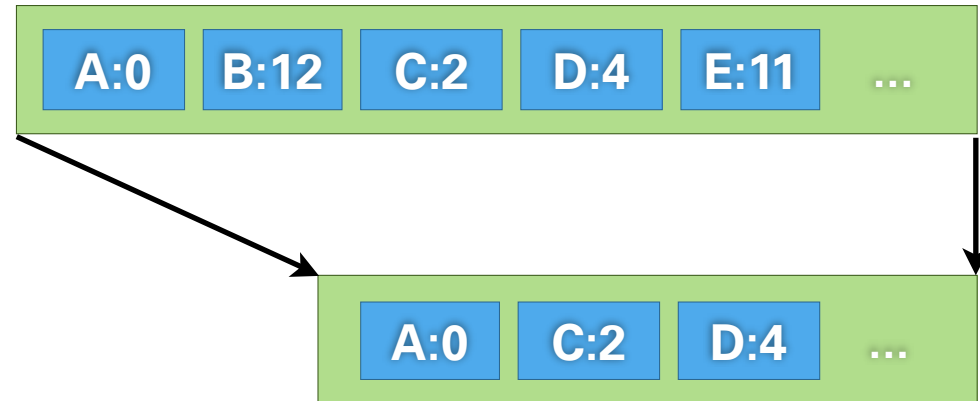
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



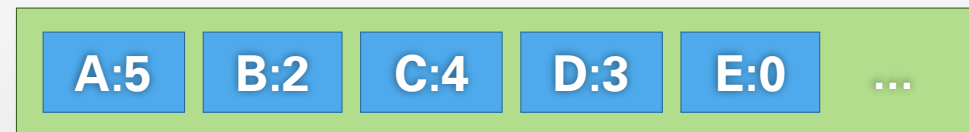
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



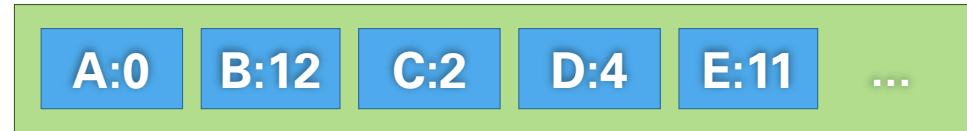
Upon receiving a window:

- Update the received entries' age (+1 for transfer)
- Update entries in local vector where newer information has been received



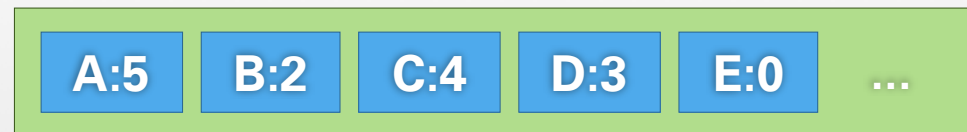
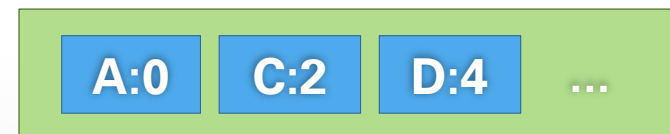
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



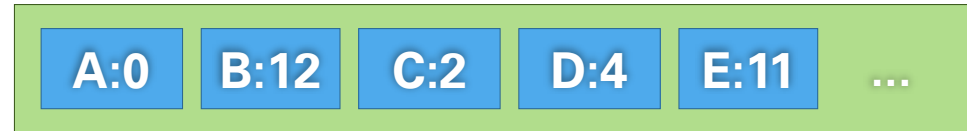
Upon receiving a window:

- Update the received entries' age (+1 for transfer)
- Update entries in local vector where newer information has been received



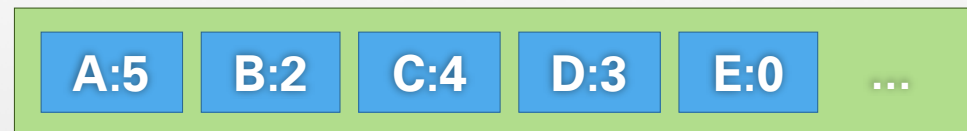
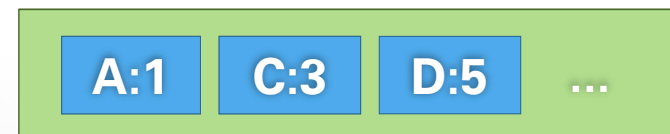
Each time unit:

- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node



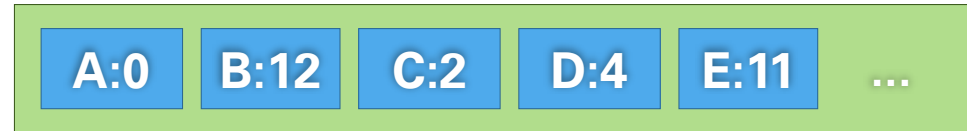
Upon receiving a window:

- Update the received entries' age (+1 for transfer)
- Update entries in local vector where newer information has been received



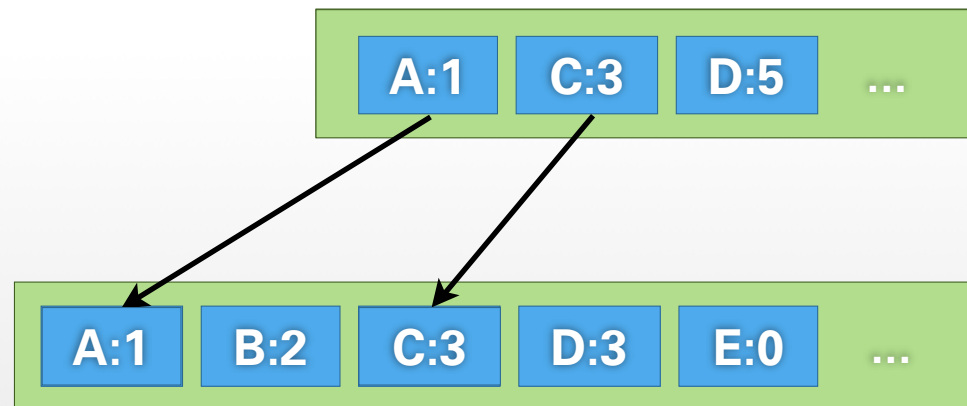
Each time unit:

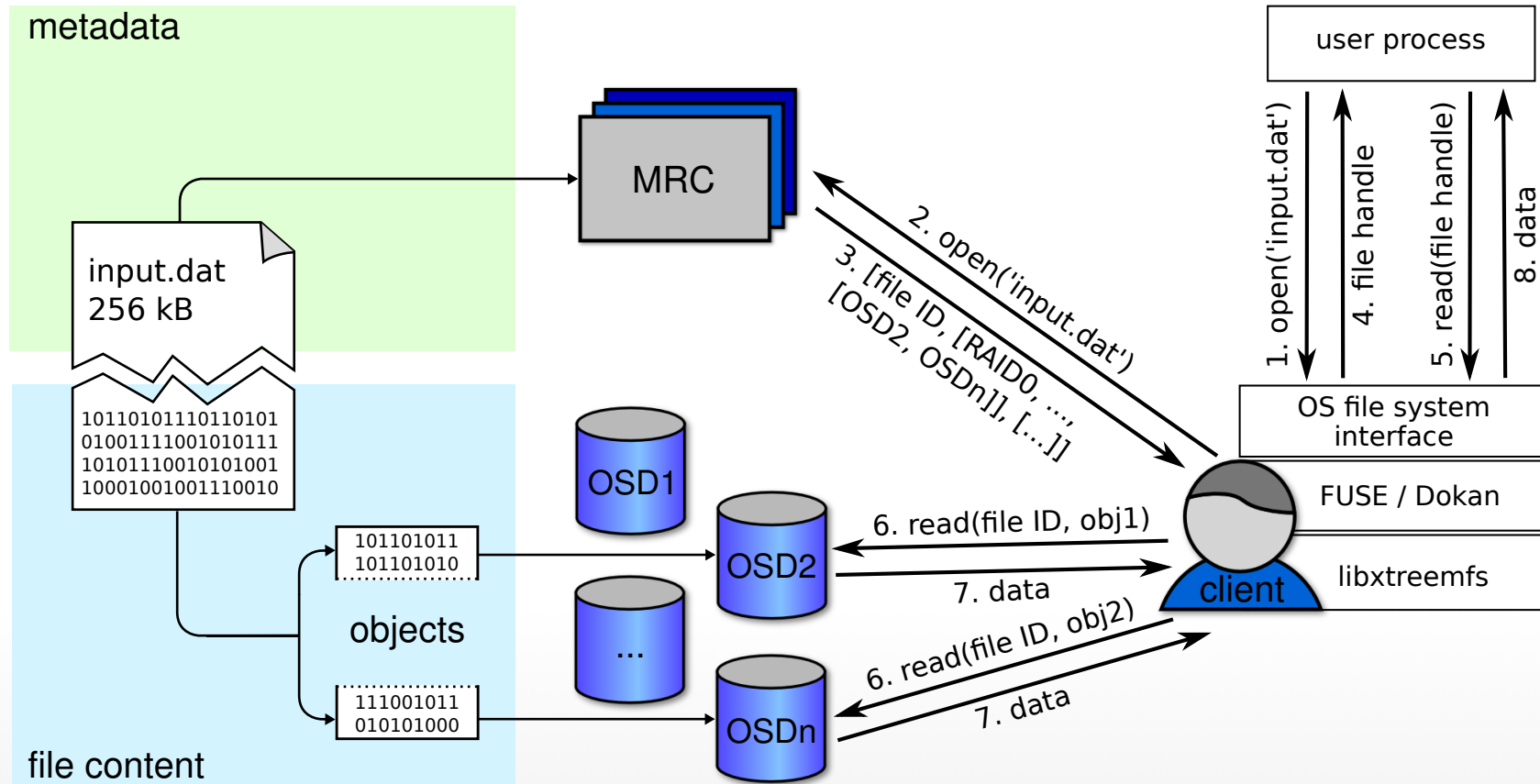
- Update local info
- Find all vector entries up to age T (called a window)
- Send window to 1 randomly selected node

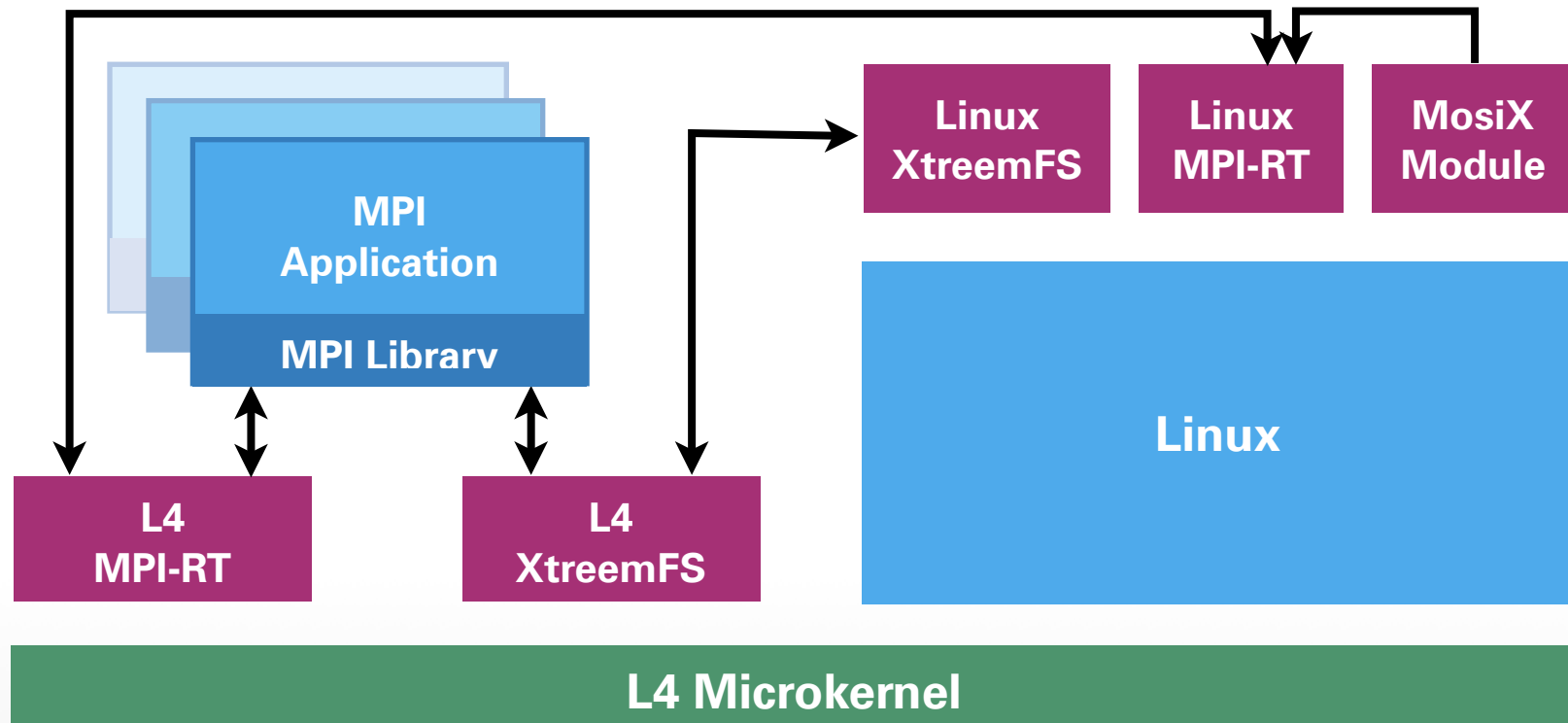


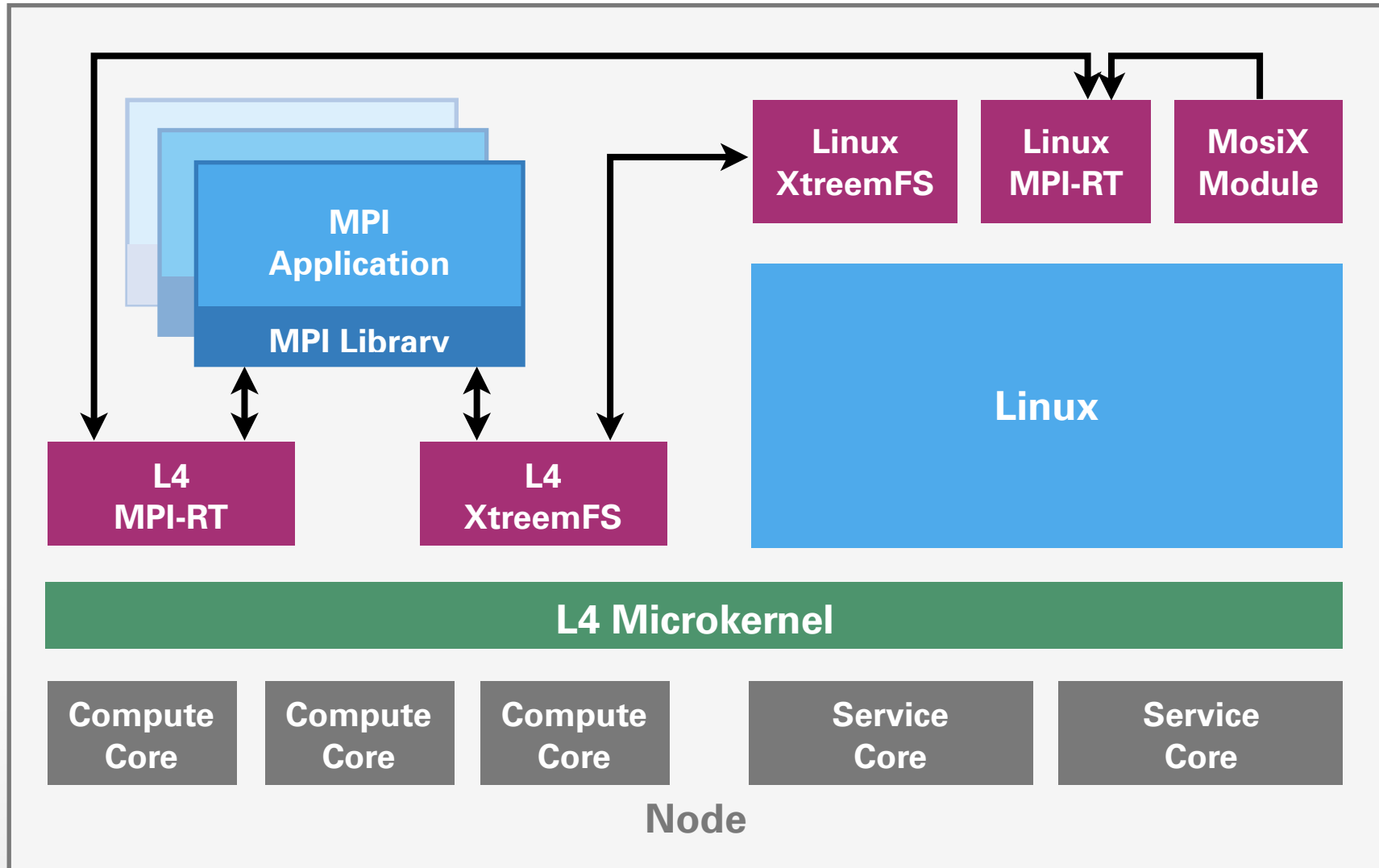
Upon receiving a window:

- Update the received entries' age (+1 for transfer)
- Update entries in local vector where newer information has been received

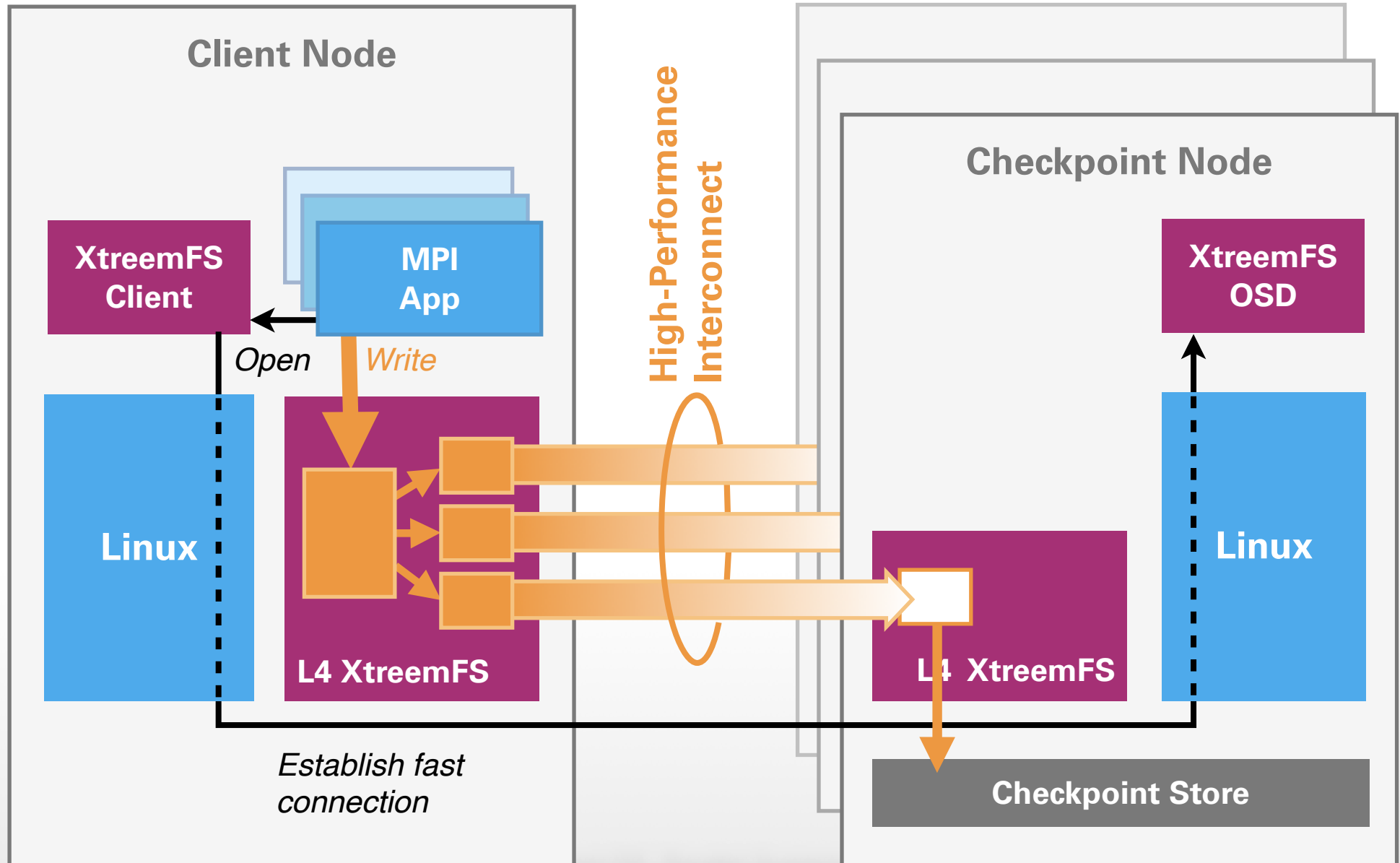


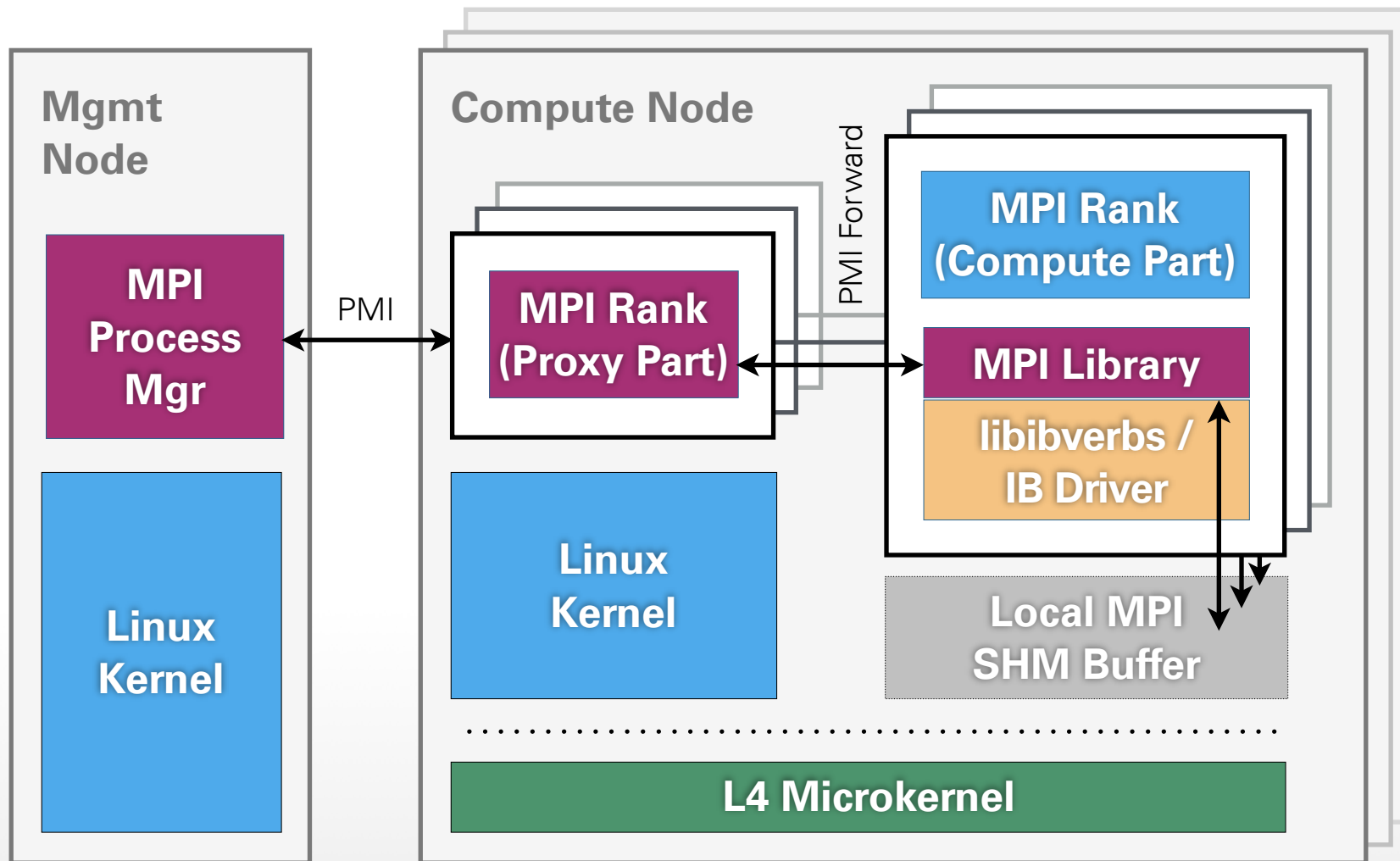


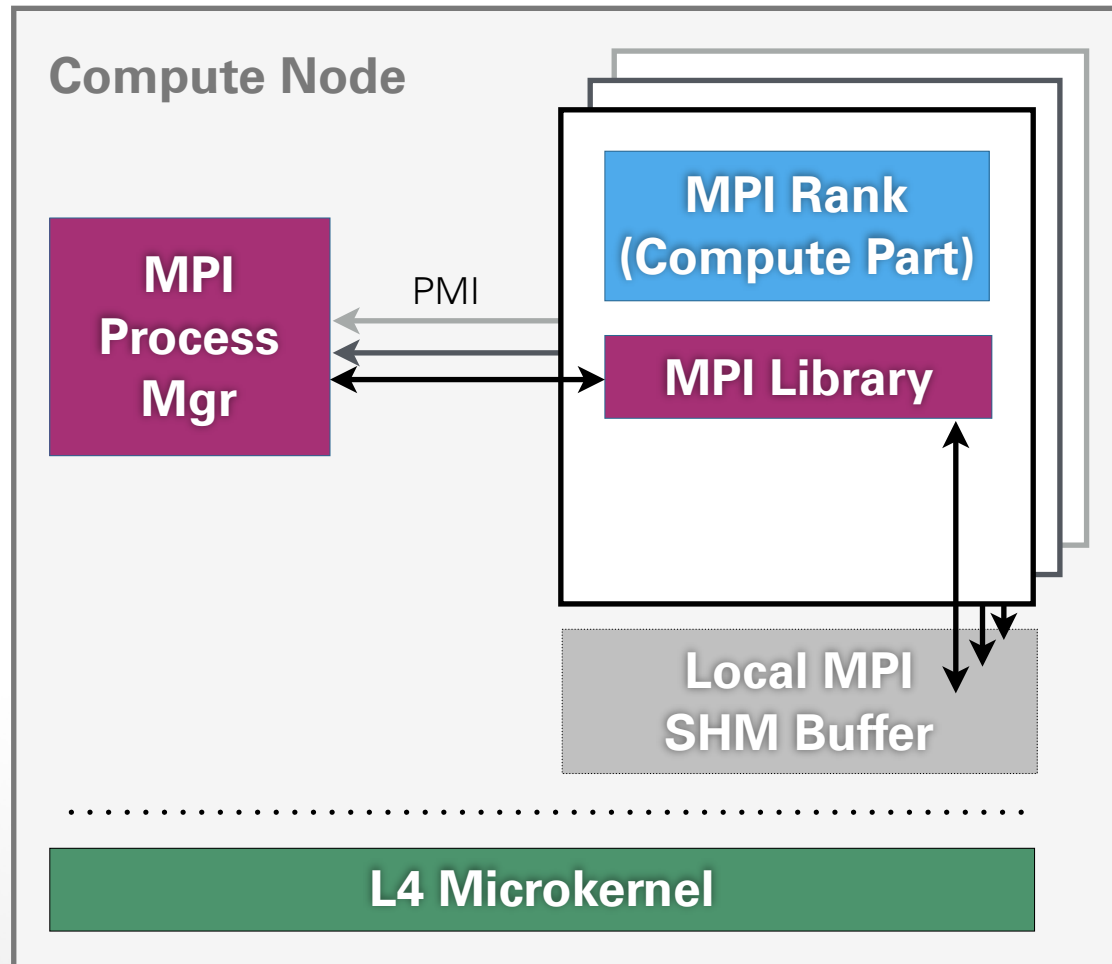




XTREEMFS: FAST PATH







DESIGN CHALLENGES

Fine-grained work splitting for system-supported load balancing?

How to synchronize? RDMA + polling ./ Block?

Gossip + Heuristics for EXASCALE ?

Application / system interface? “Yell” for help?

Compute processes, how and where to migrate / reroute communication?

Replication instead of / in addition to checkpoint/restart?

Reuse Linux (device drivers)?

- Perf counters for Network
- fast redirection of messages
- flash on node circumventing FTL
- quick activation of threads without polling