

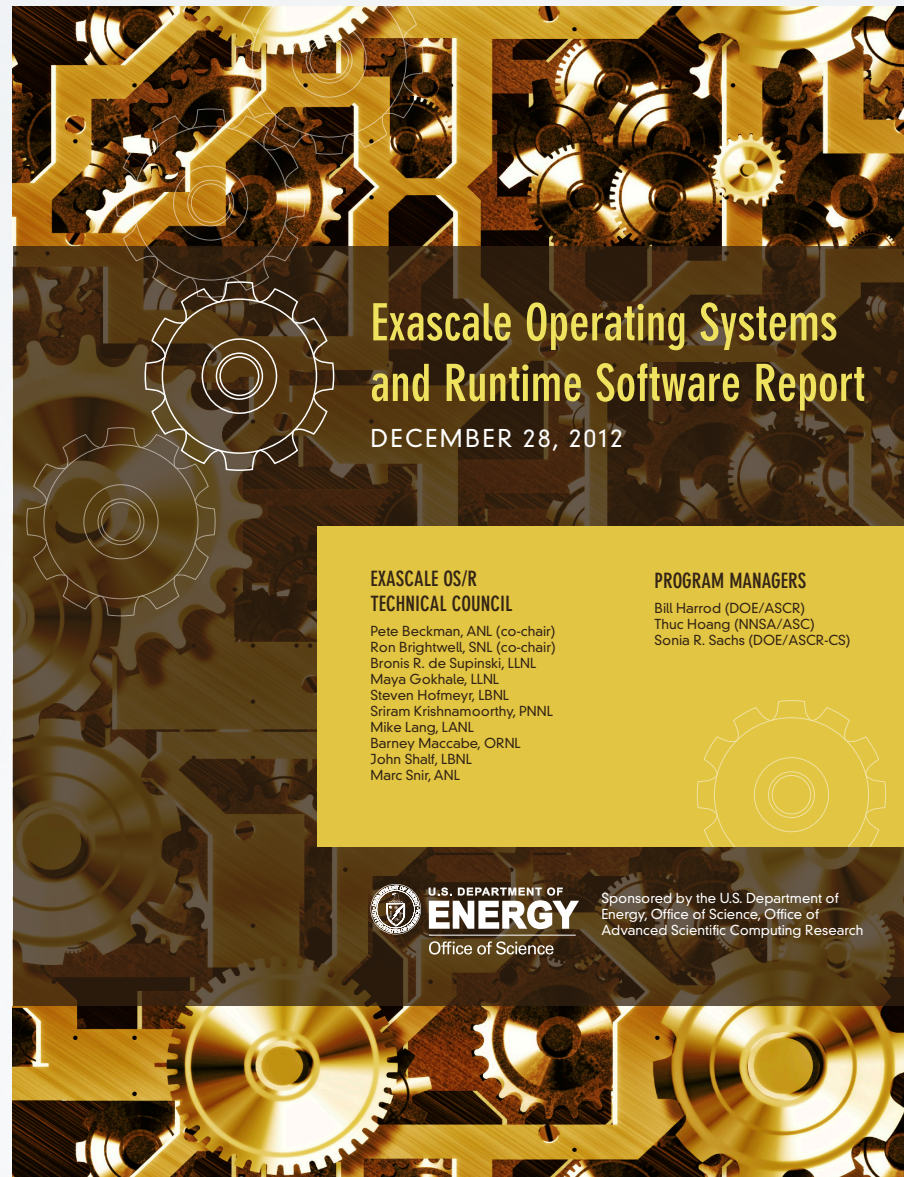


**TECHNISCHE
UNIVERSITÄT
DRESDEN**

Faculty of Computer Science Institute of Systems Architecture, Operating Systems Group

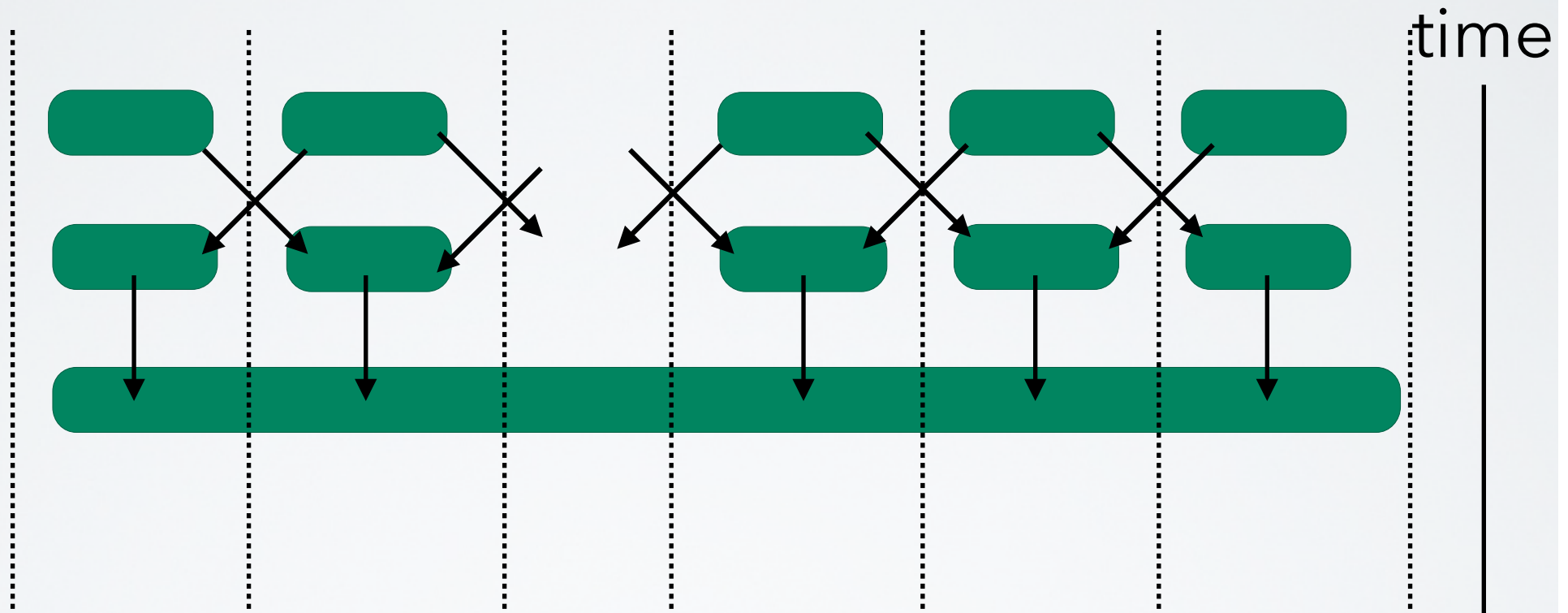
BUILDING BLOCKS FOR AN EXA-SCALE OPERATING SYSTEM

**HERMANN HÄRTIG
ROSS 2014**



the ideal world assumption

- identical predictable reliable nodes
- fast deterministic reliable interconnect with isolated partitions of fixed size
- balanced applications



applications split into fixed-size chunks of work
one thread/core

systems software:

optimize communication latency

- RDMA & busy waiting
- batch scheduler for start / stop
- separate servers for IO
small OS on each node

no OS on critical path

observations (infiniband cluster):

- CPUs run at different speeds:
3 of 16 @ 1.2GH, others 13 @ 2.9GH
- “turbo boost” switched off
- measurements not reproducible
(node allocations arbitrary)
- application restart due to failures

Application: COSMO-SPECS+FD4



Unbalanced
compute times of
ranks per time step



Hand balanced
compute times of
ranks per time step

- very large number of processing units
- diverse on-node storage
- heterogeneous processing units
- different and dynamic speed
- higher failure rate
- dark silicon
- new resource trade-offs
- dynamic application behavior

Fast and Fault-Tolerant

Microkernel-Based Operating System



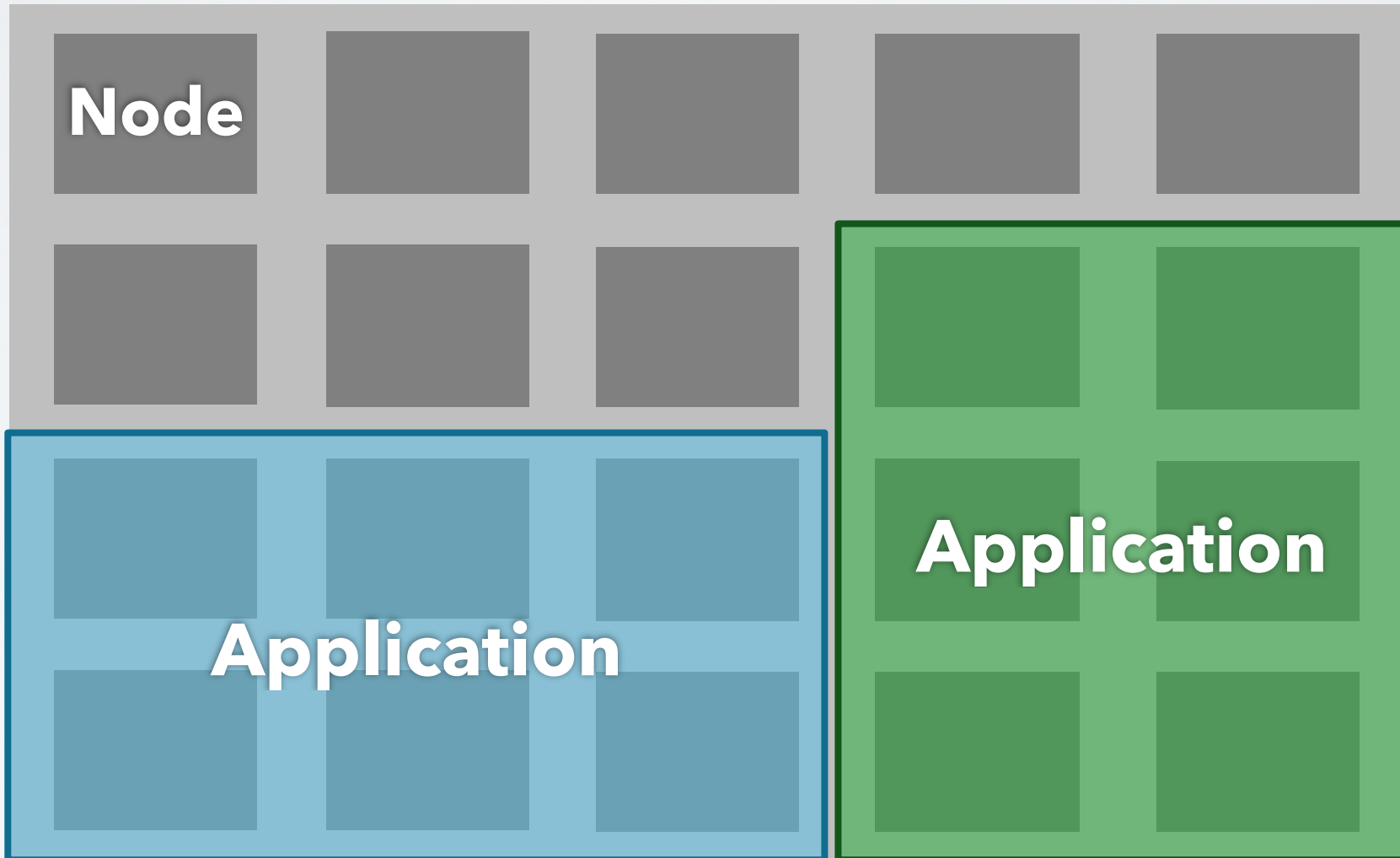
The Hebrew University
of Jerusalem

DFG Deutsche
Forschungsgemeinschaft

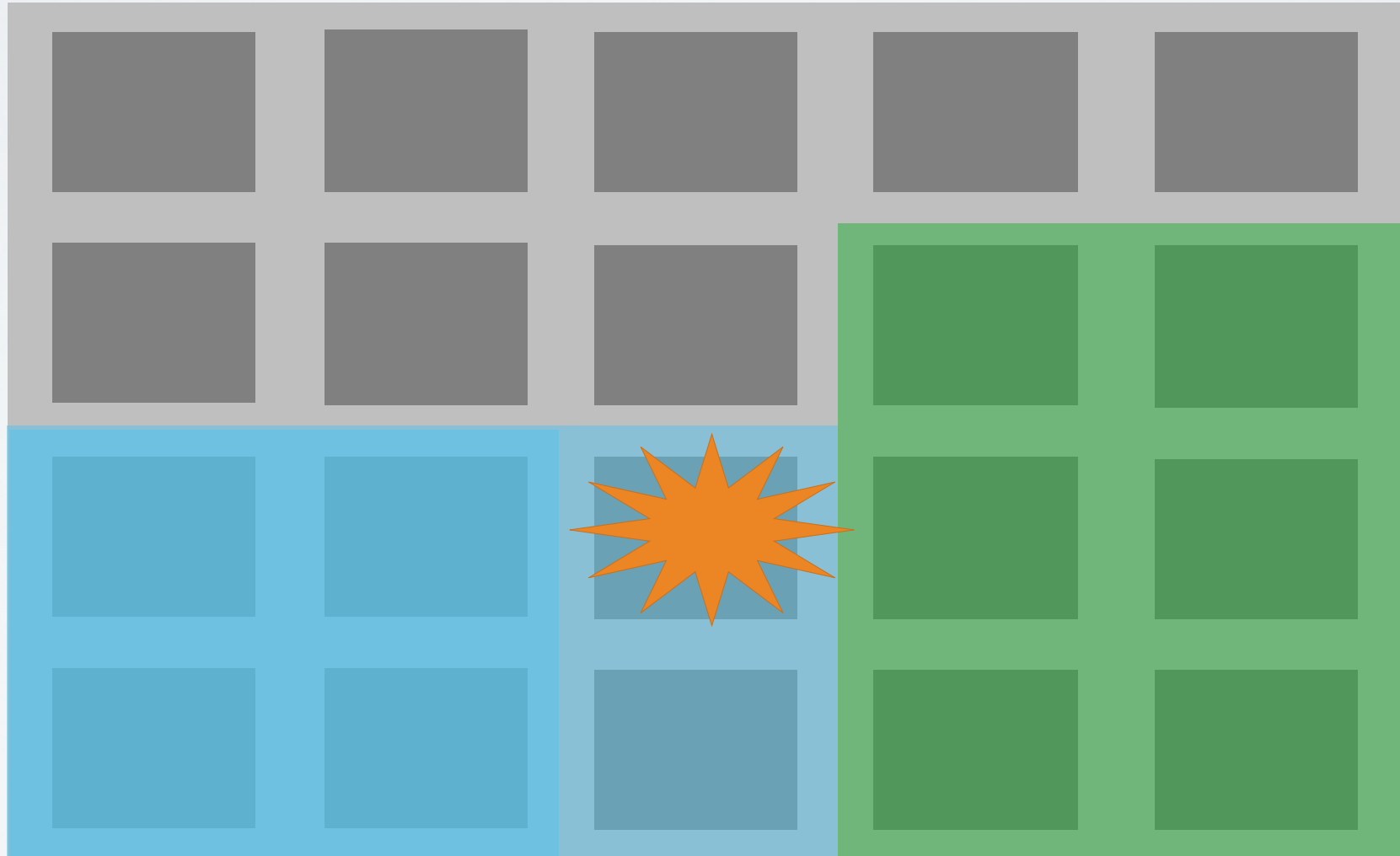


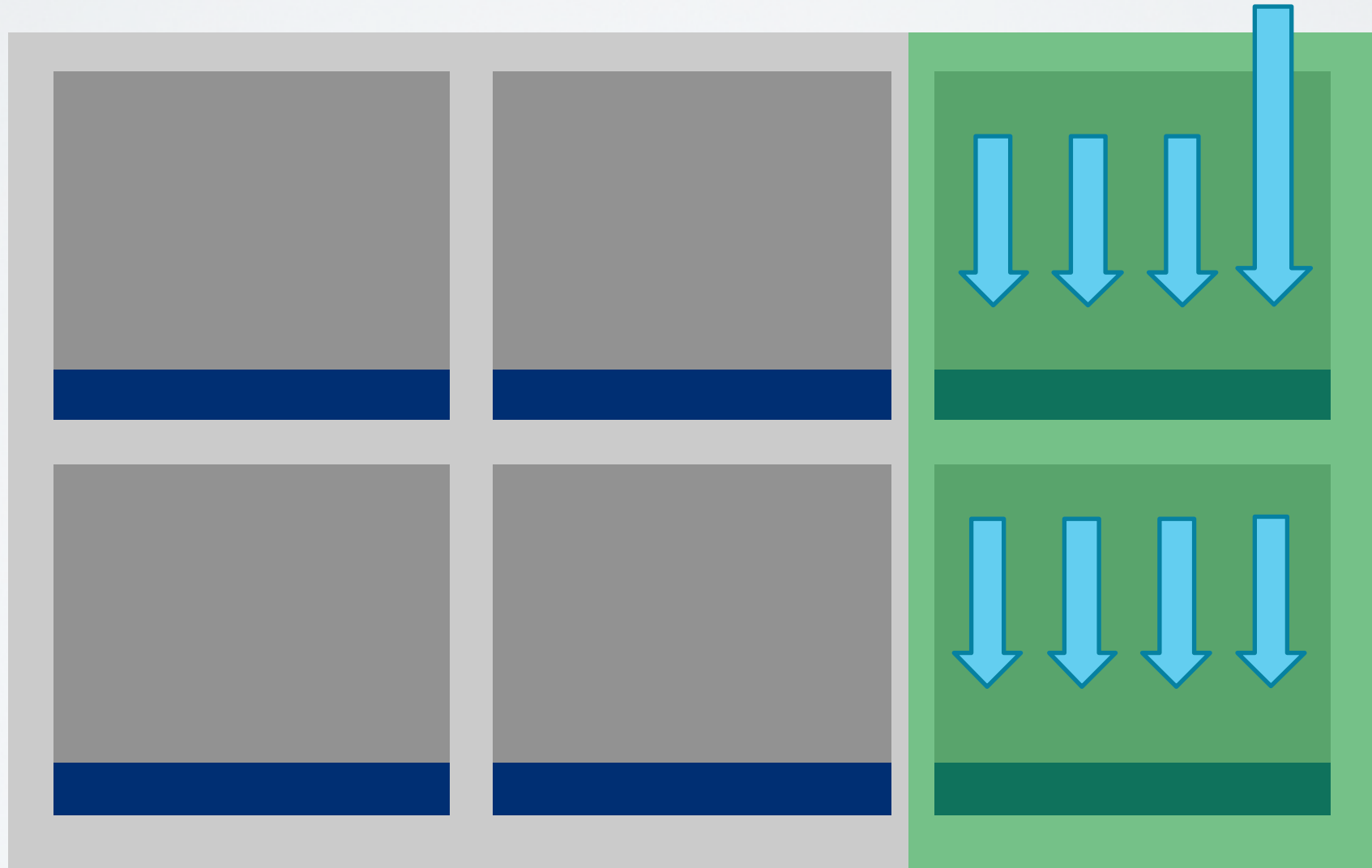
German Priority Programme 1648
Software for Exascale Computing

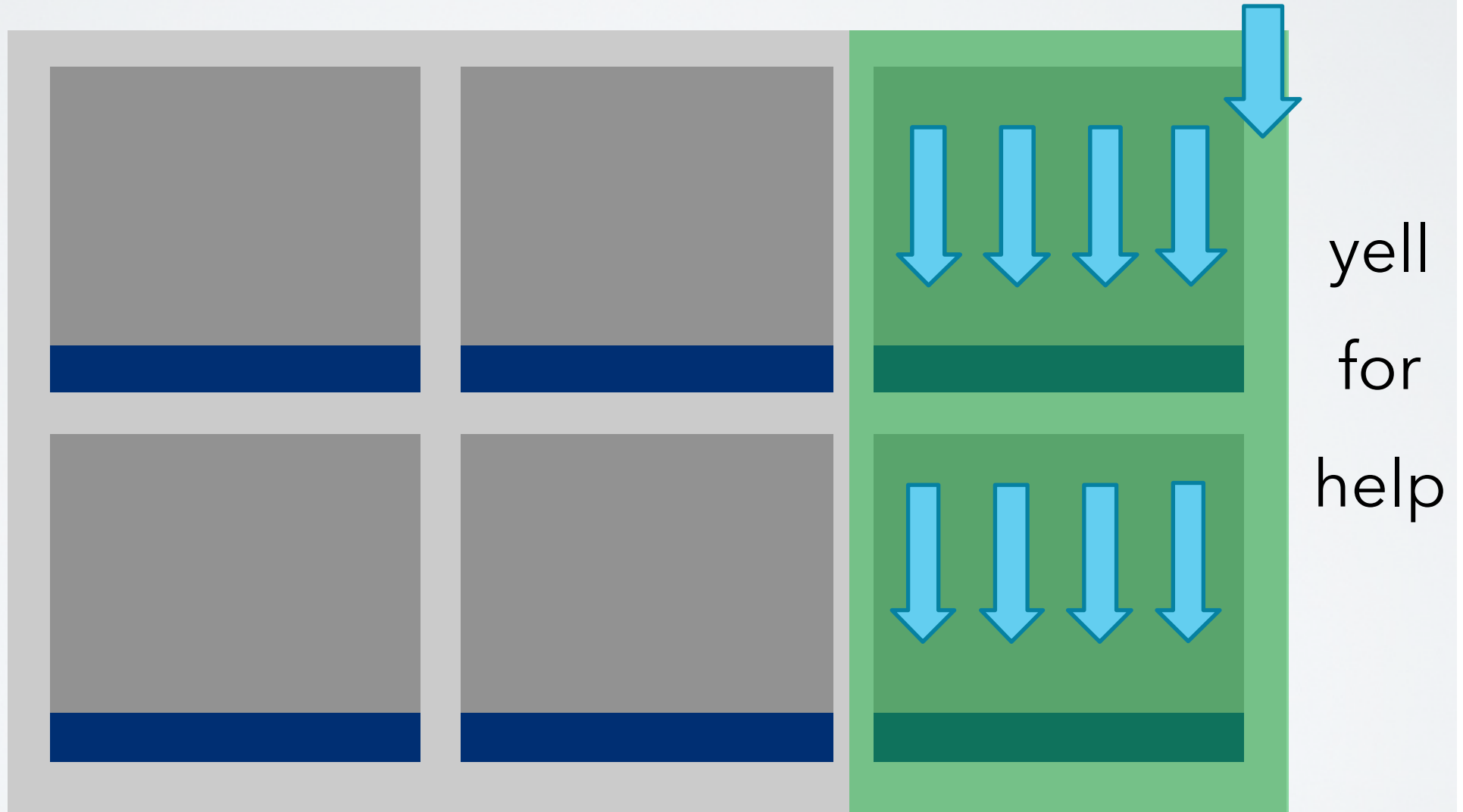


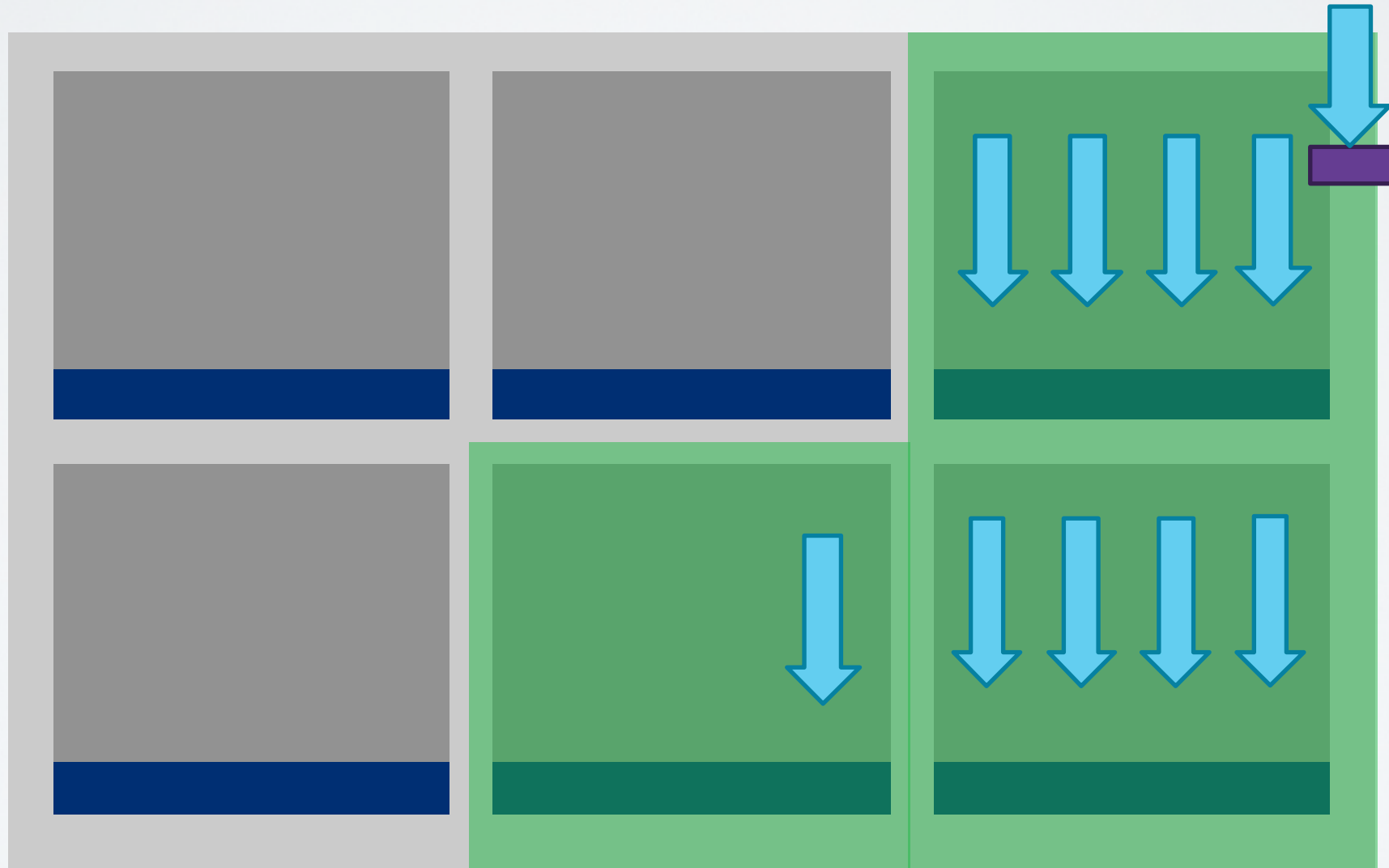


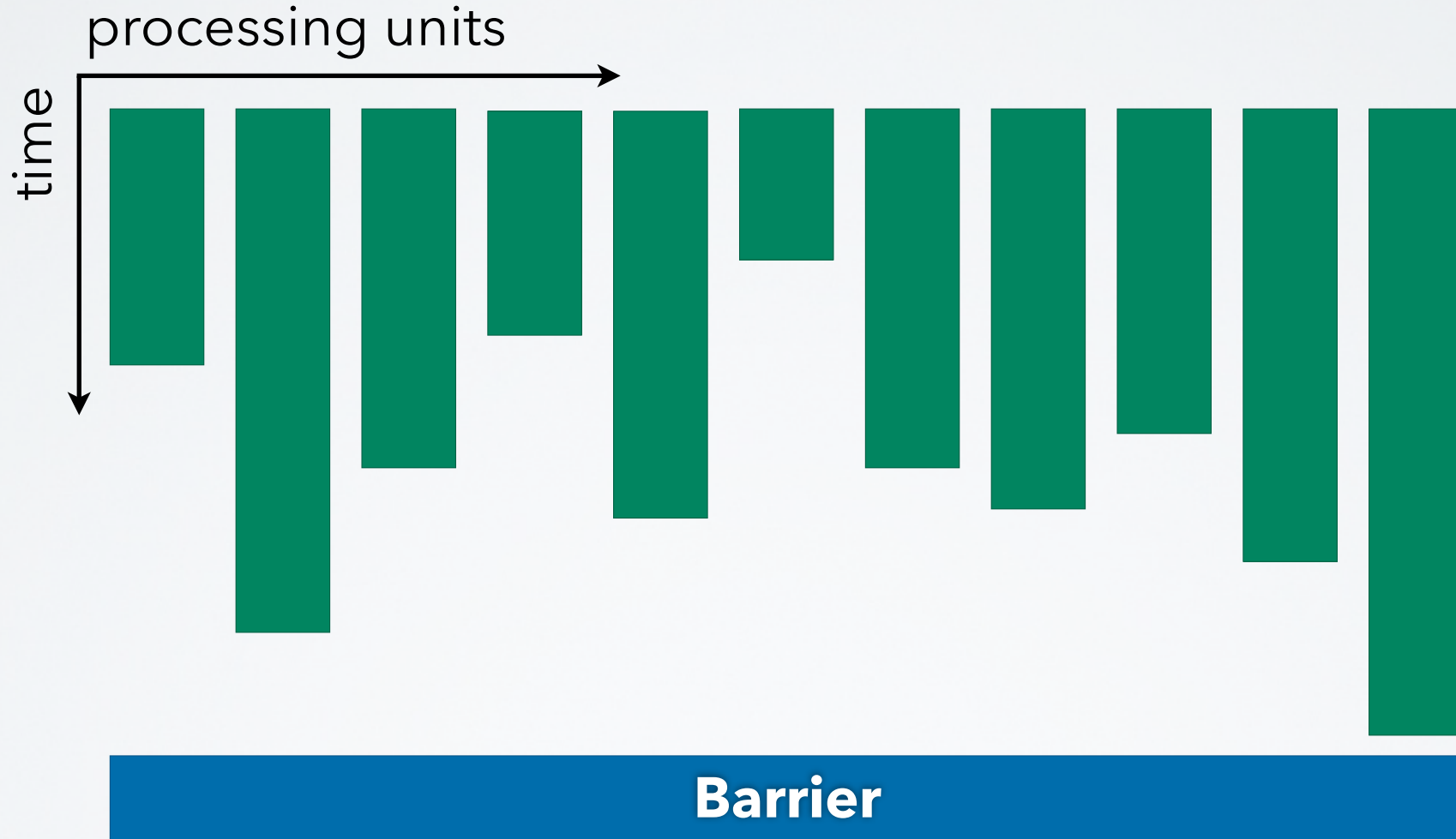


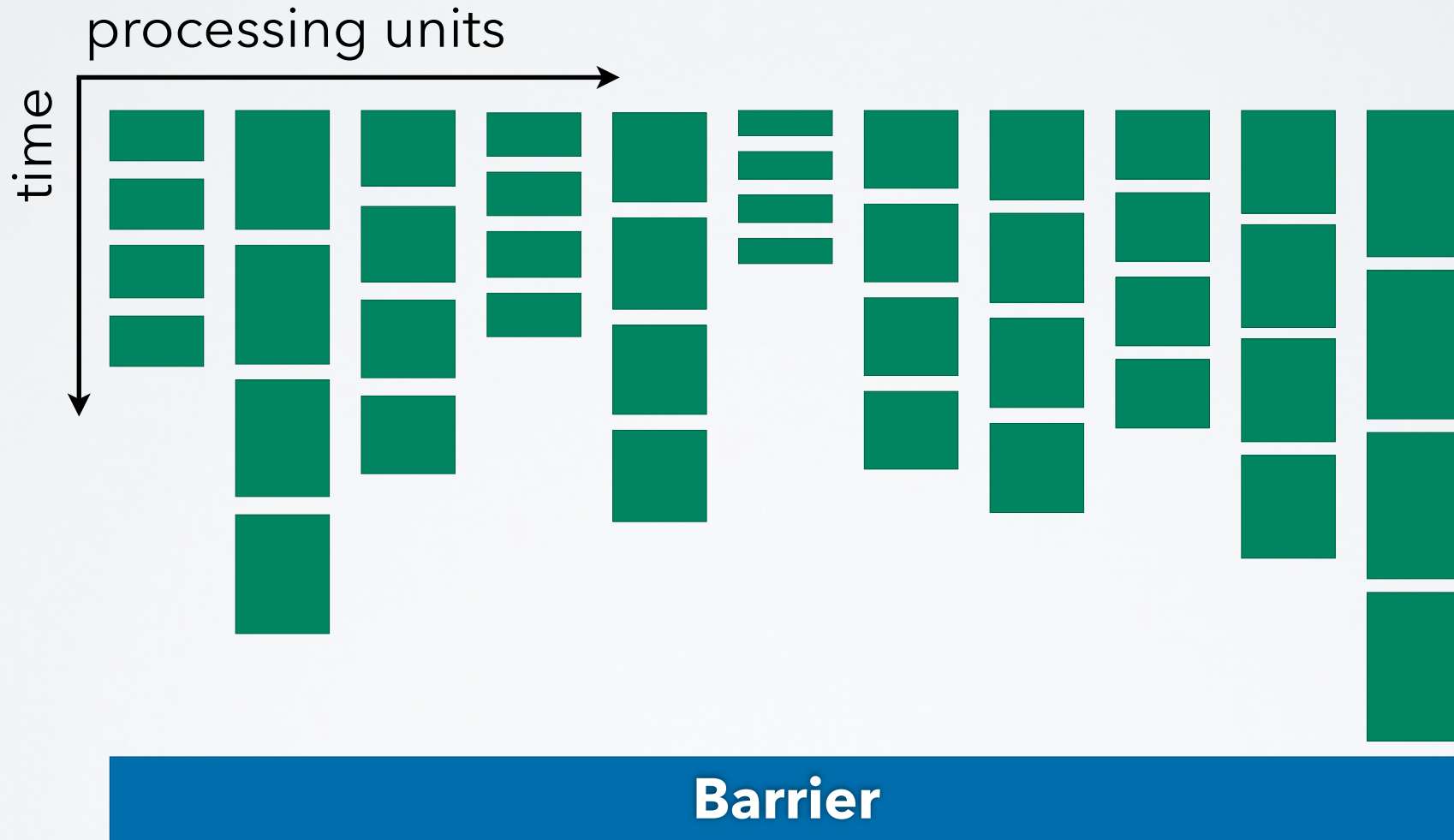


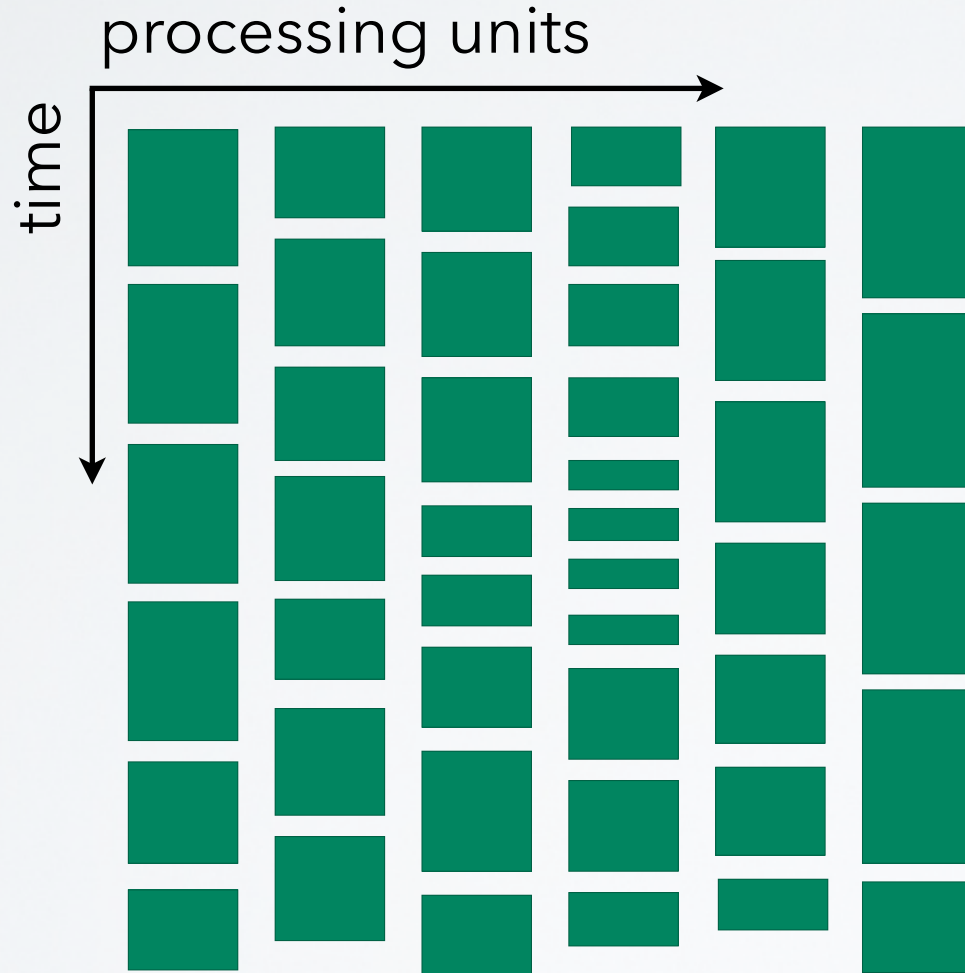


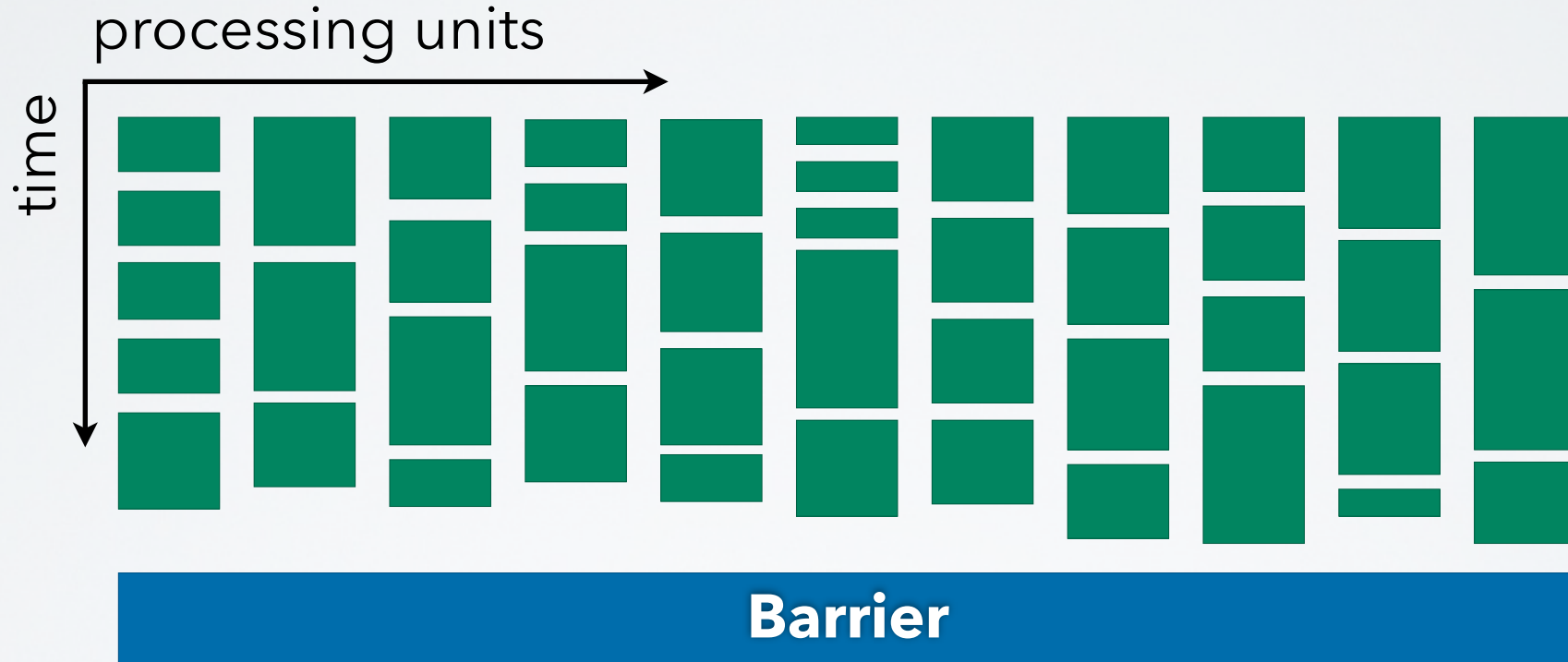


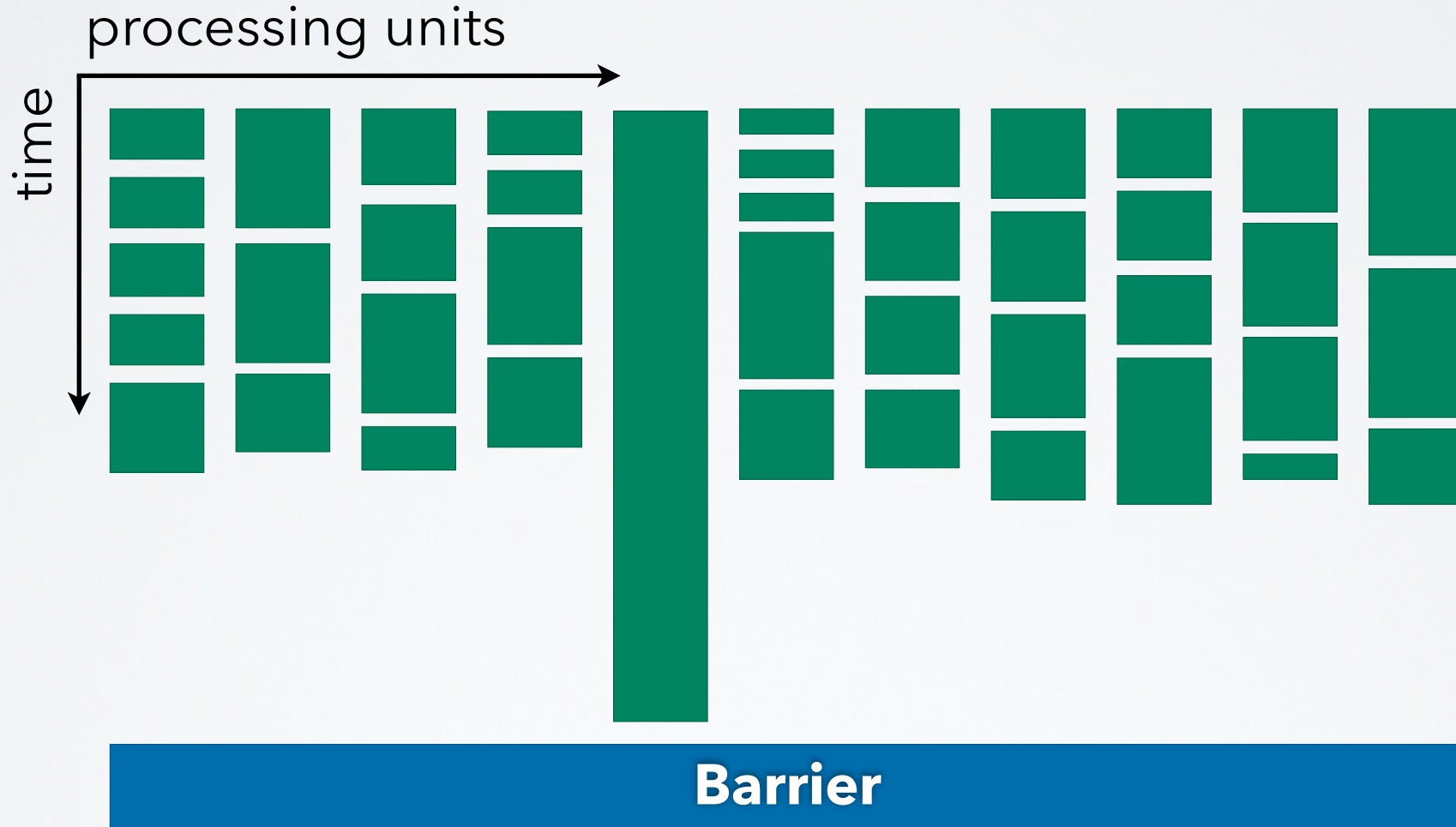


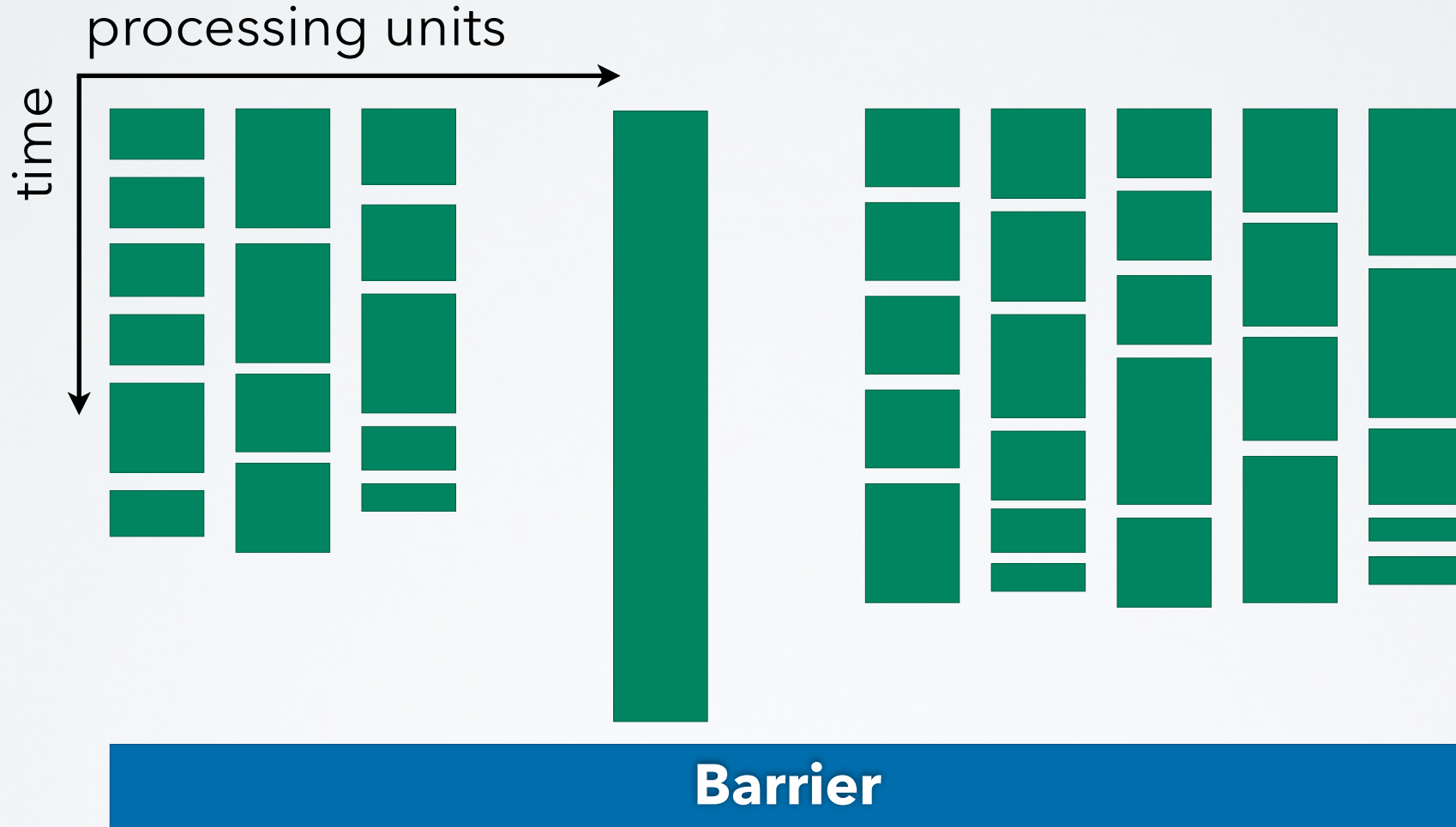








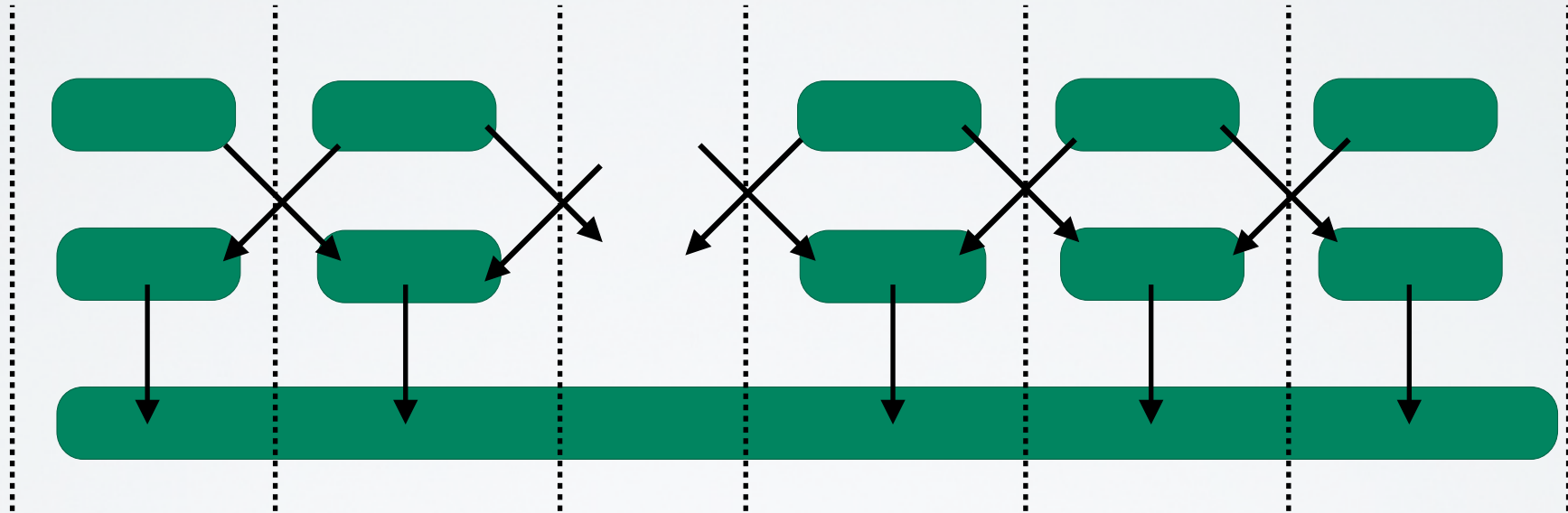




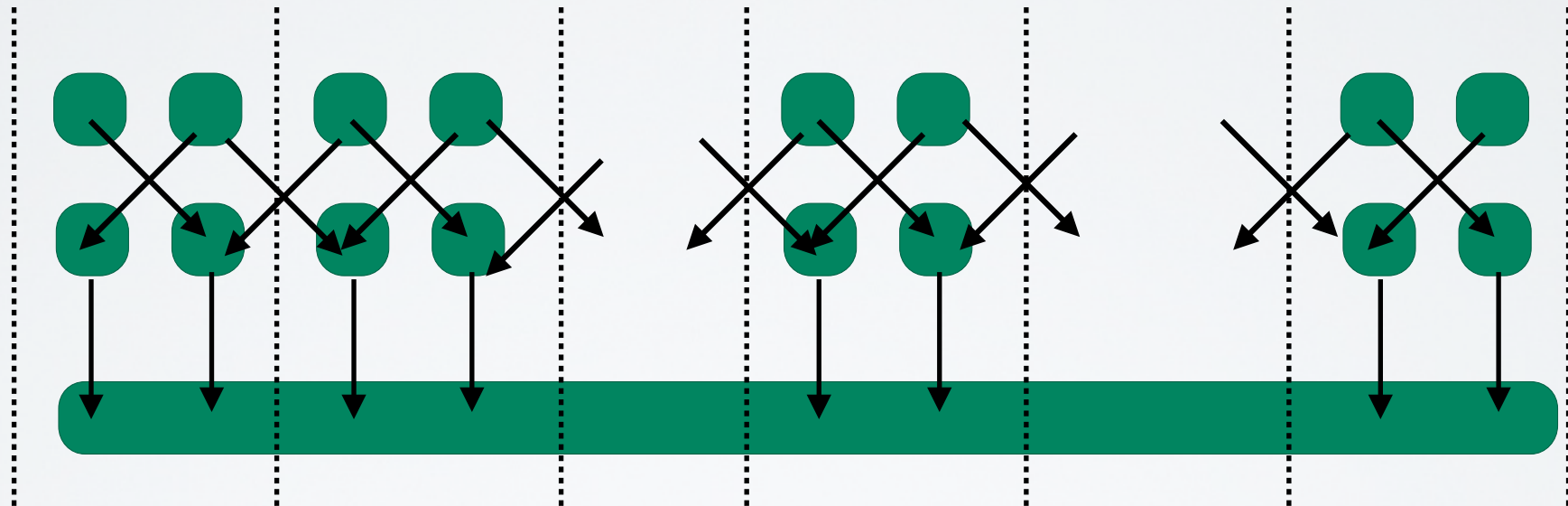
simplistic experiment

(Linux on IB cluster)

- split application into more chunks than cores
- run more MPI processes than cores

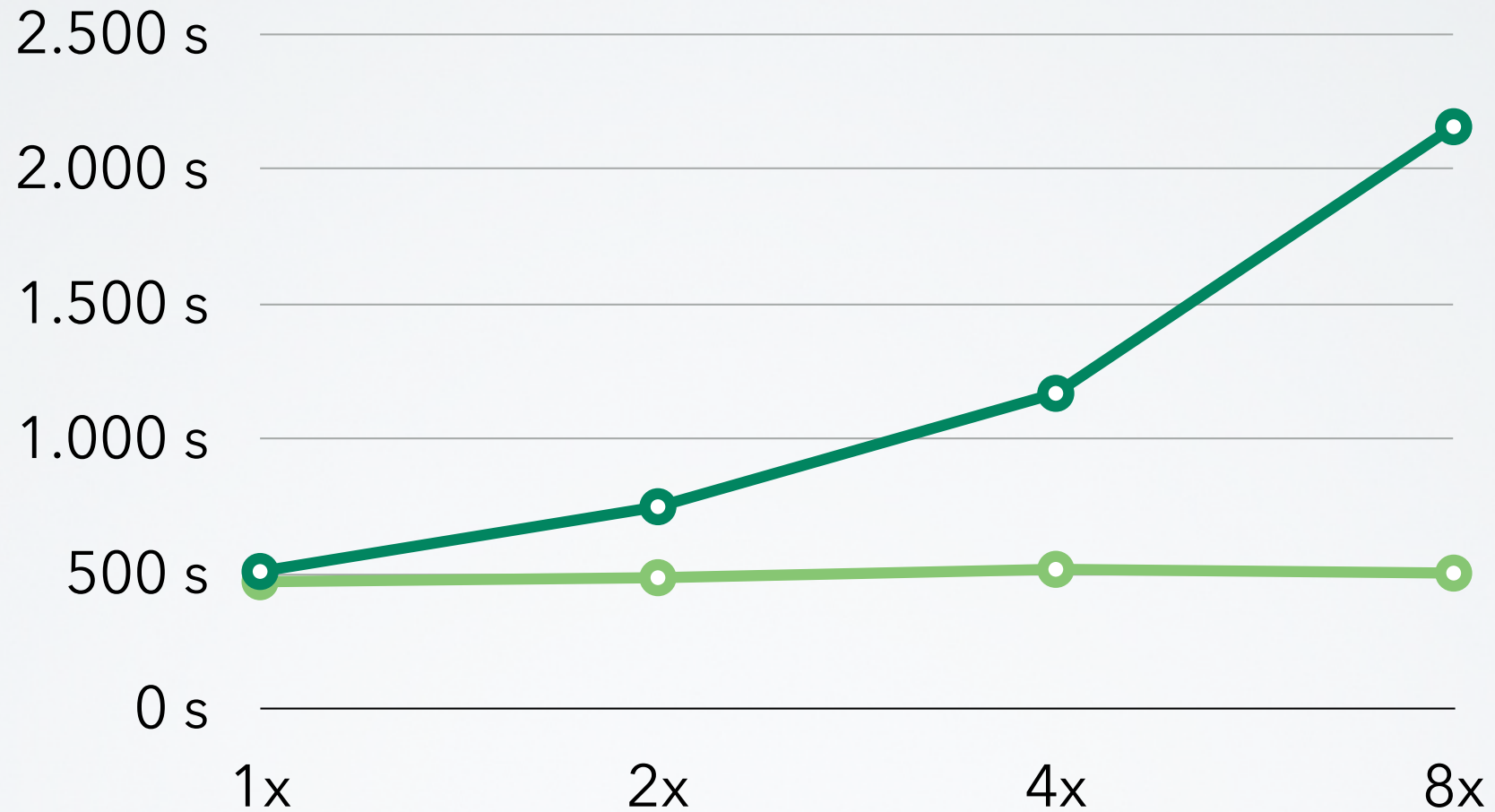


1 MPI process per core

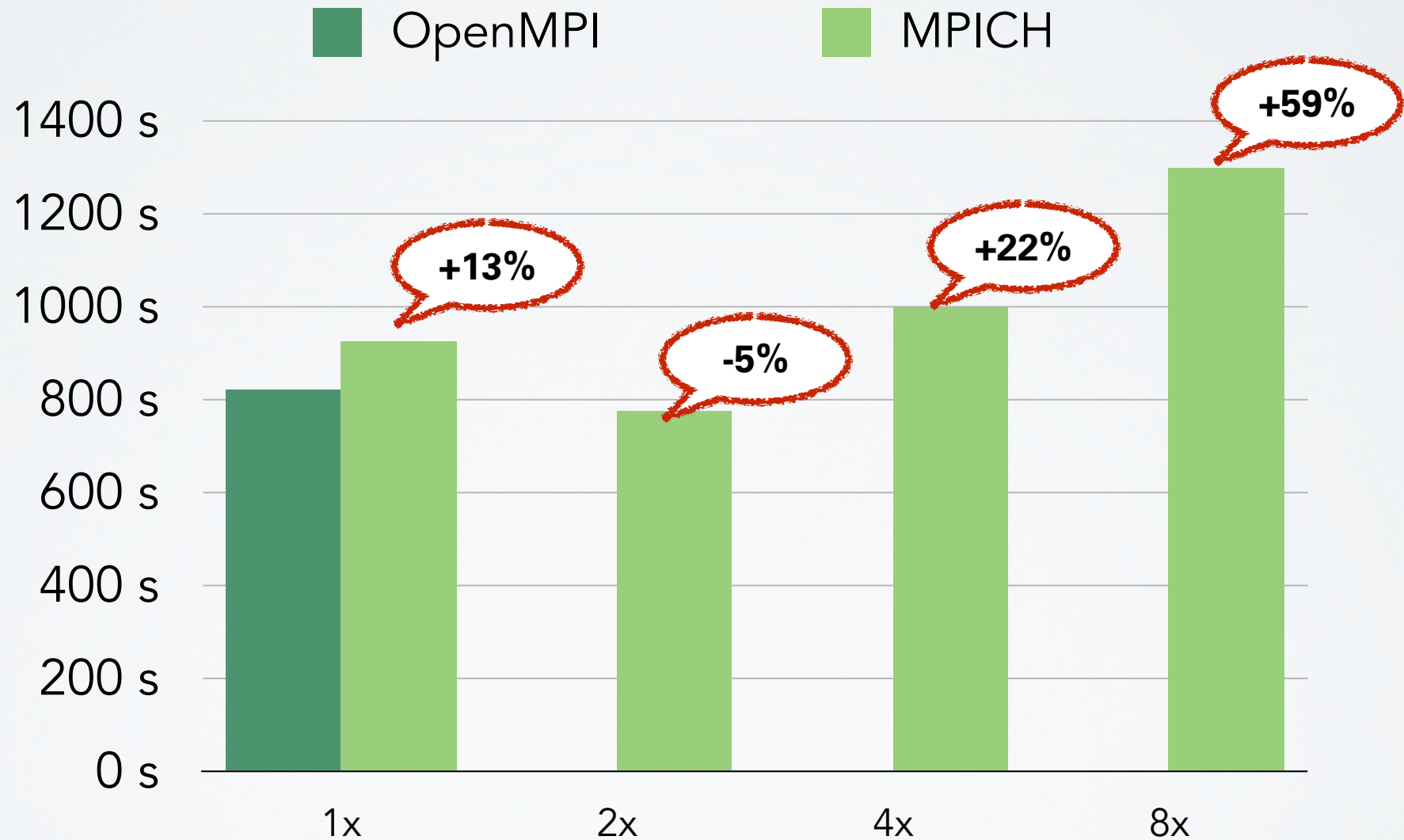


Application: COSMO-SPECS+FD4 (no load balancing)

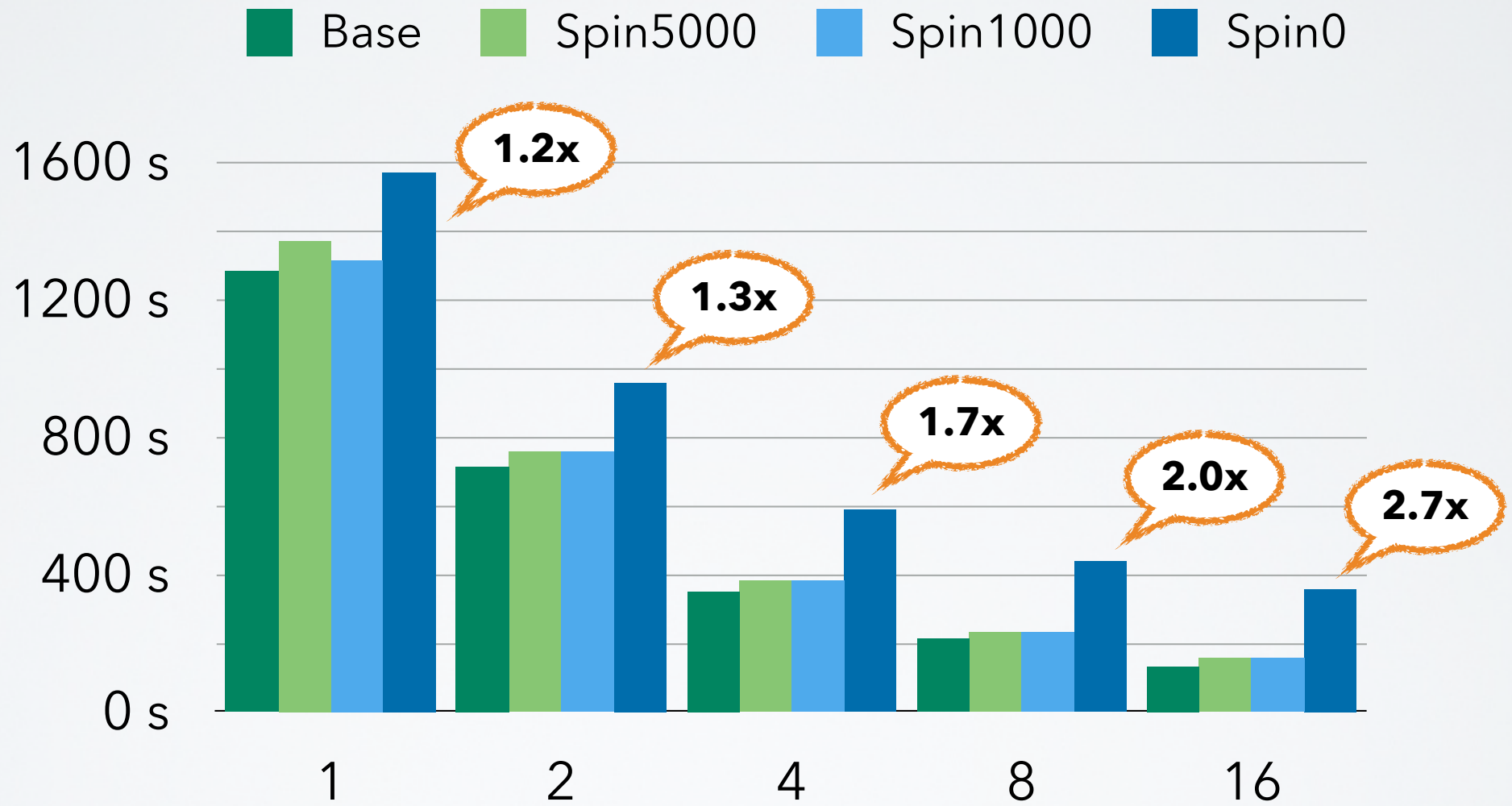
- IB-Cluster 16 nodes w/ 16 Xeon E5-2690 (Sandy Bridge) @ 2.90GHz
- 1x - 8x oversubscription (same problem size)



Oversubscription factor (more ranks)



COST OF BLOCKING



GROMACS on MVAPICH

Number of nodes

Table 2 Receive 1 Gbps latency breakdown

Description of receive packet activities	Source	Time (ns)
MAC filter determines target packet is for this machine	Estimation	200
NIC starts DMA packet header and payload into memory	Estimation	400
NIC interrupts core with MSI-X packet to APIC	Estimation	500
Hardware MSI-X interrupt service routine to parse what caused interrupt	Estimation	270
Interrupt cause register read requirement	Measurement	1,000
ISR packet processing of descriptor to update receive queue	Measurement	300
SoftIRQ (deferred procedure call in Windows)	Measurement	1,287
TCP and IP receive side processing	Measurement	570
Wakeup application to process socket information	Measurement	1,274
Kernel to application space data copy	Measurement	208
ACK the pong received by the remote sender	Measurement	1,117
Application receive message overhead to register completion	Measurement	621
Total receive packet time		7,747

Architectural Breakdown of End-to-End Latency in a TCP/IP Network

Int J Parallel Prog, 2009 (intel)

overdecomposition overhead:

- blocking overhead
- additional messages

fast inter process communication

“Lightweight message and thread management”



L4 MICROKERNEL FAMILY



Jochen Liedtke
1953 -2001



Simko3 aka "Merkel phone"



Franklin eBookMan

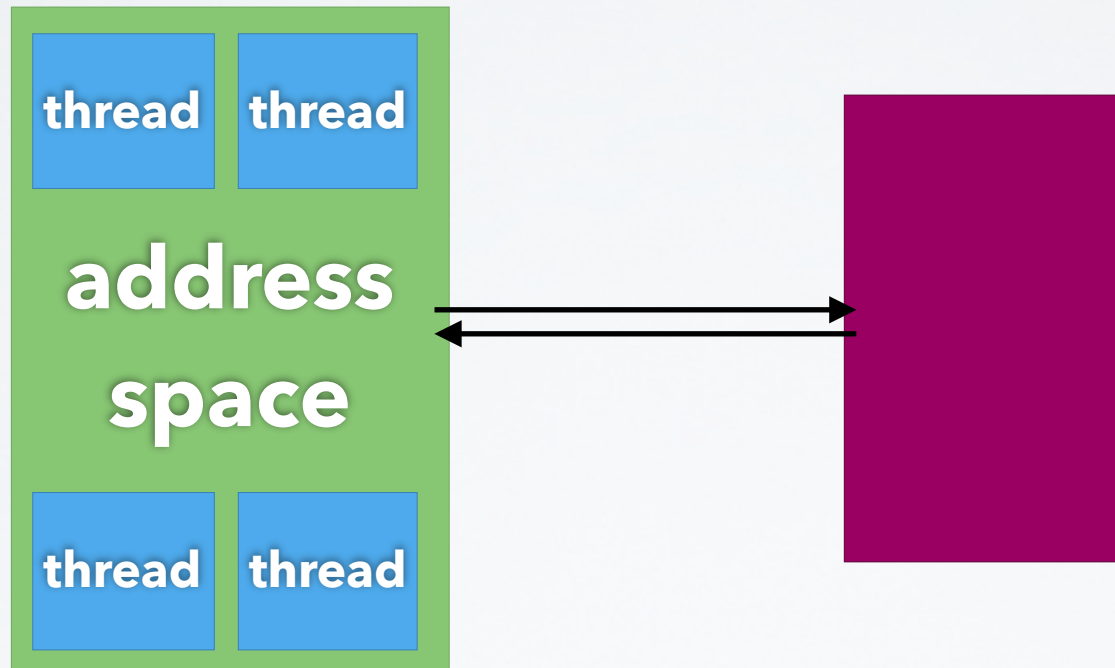
nearly 2 billion
Qualcomm boards

Cellphone Baseband processors

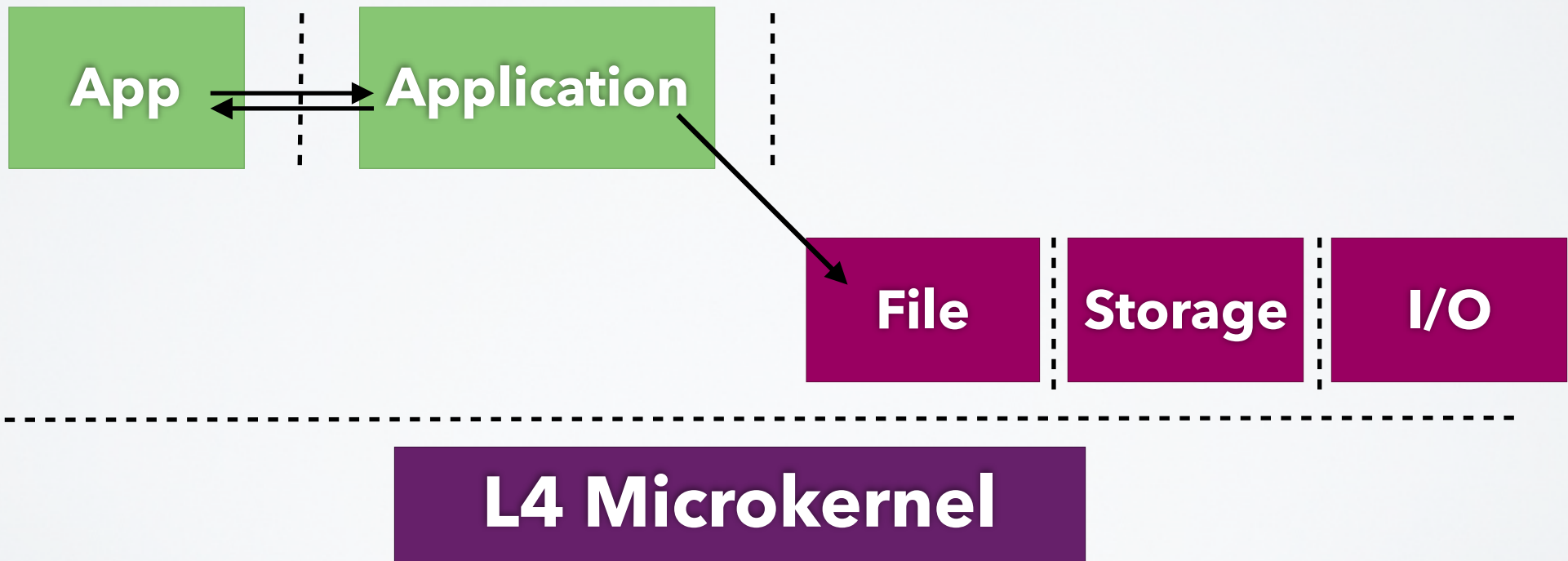


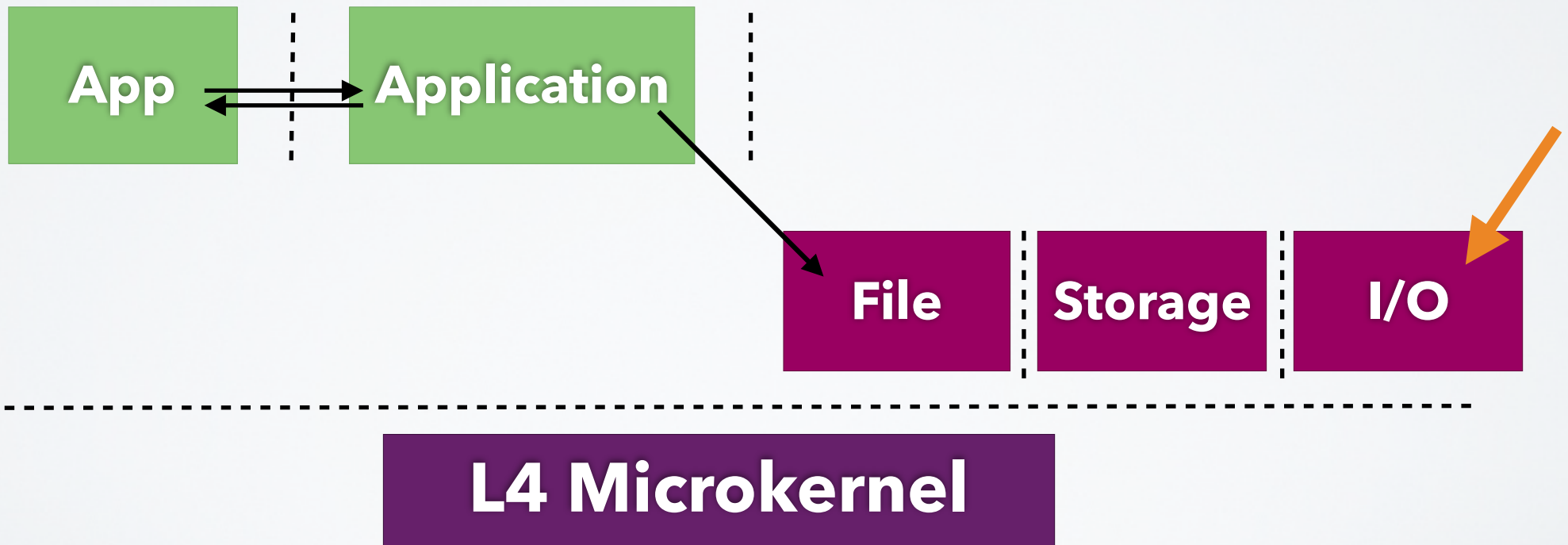
Airbus

3 ABSTRACTIONS



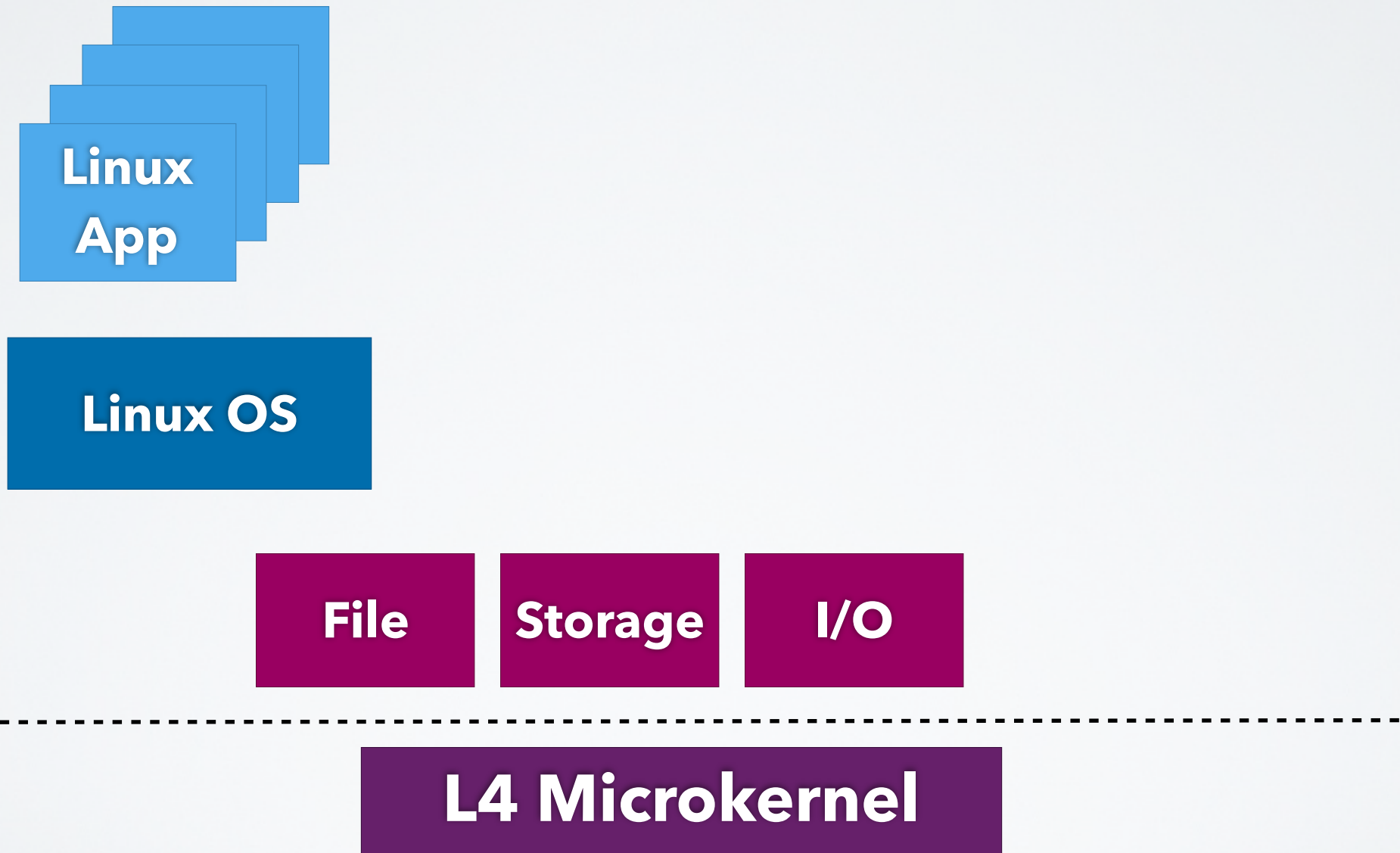
L4 Microkernel





interrupt response (unblocking):

- HW-IRQ → Kernel-Handler → SemUp → iret idle → schedule → iret user_handler
- 900 cycles, 0.3 micro second (best case)
4500 cycles (in deeper sleep level)
- Quad-Core Intel(R) Core(TM) i7-4770
CPU @ 3.40GHz

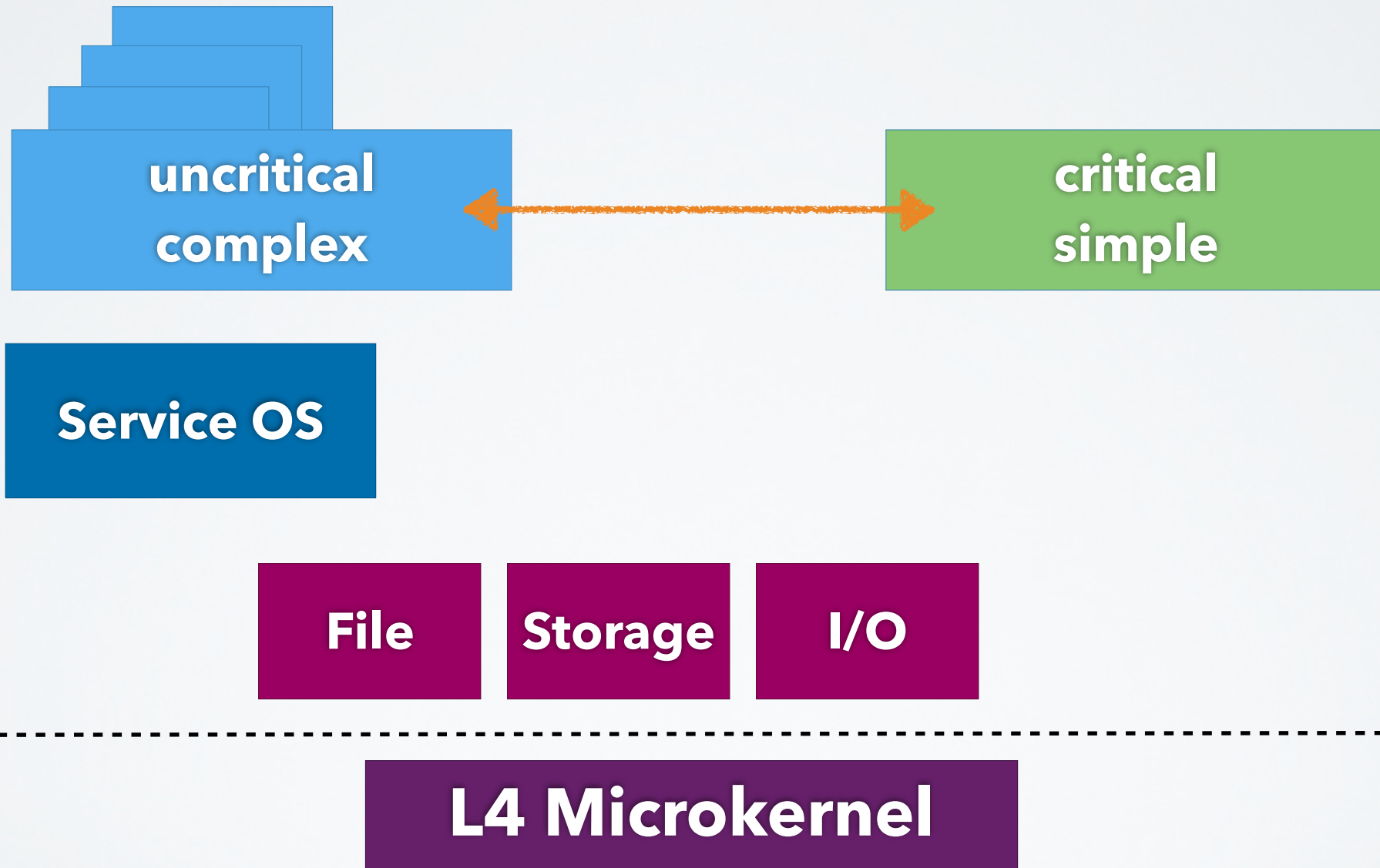


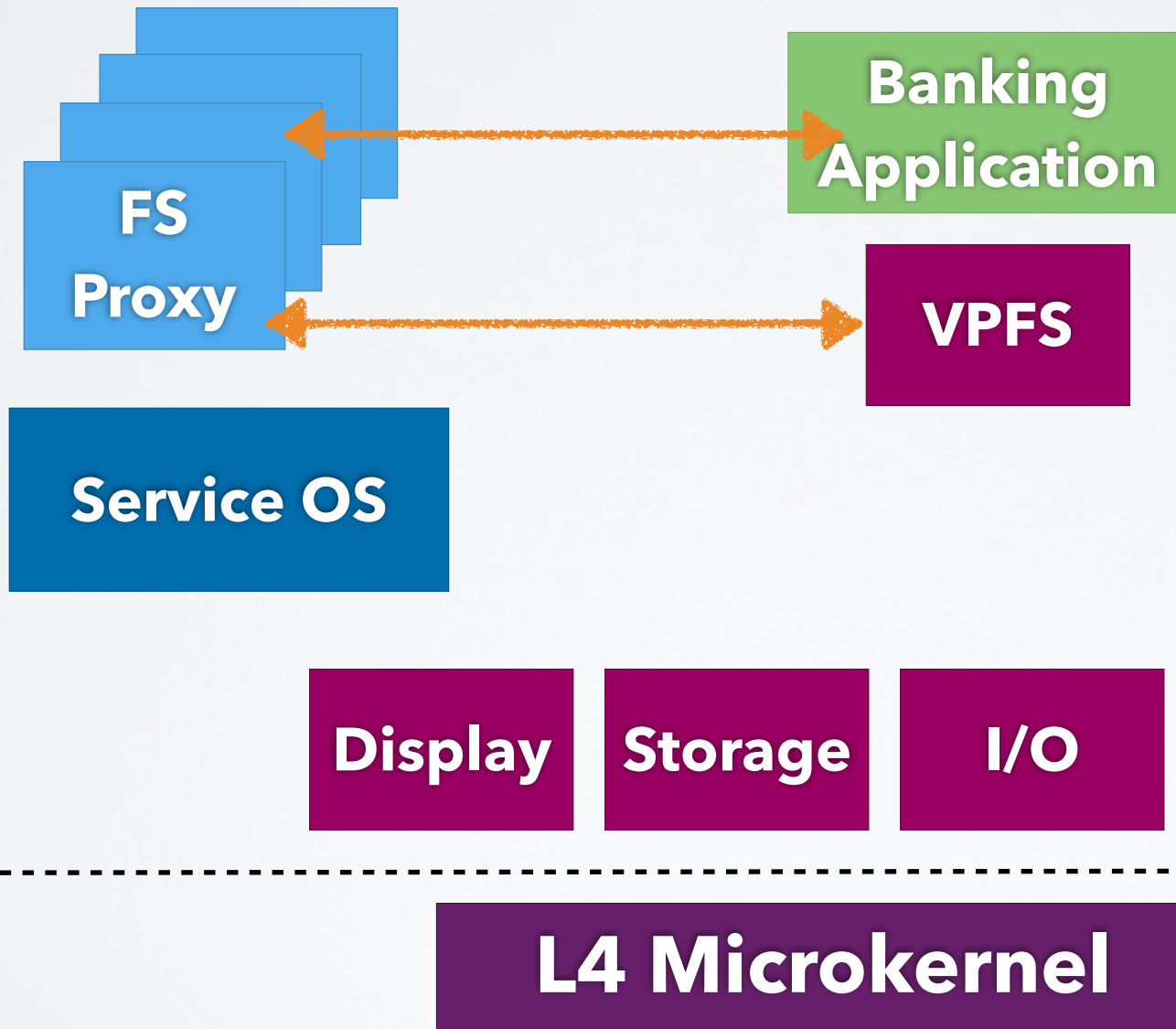


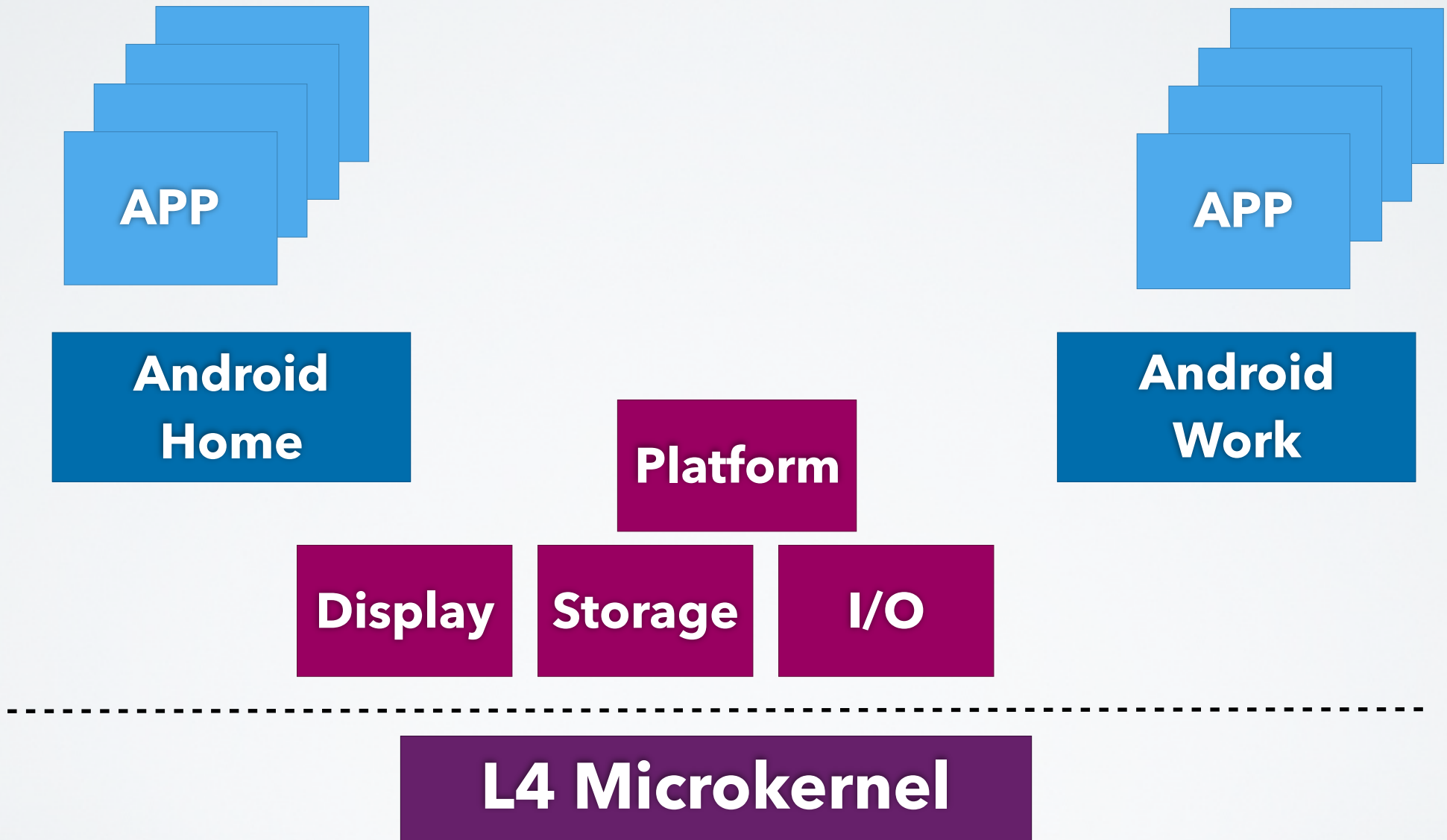
real-time

security: small Trusted Computing Base

resilience: small Reliable Computing Base

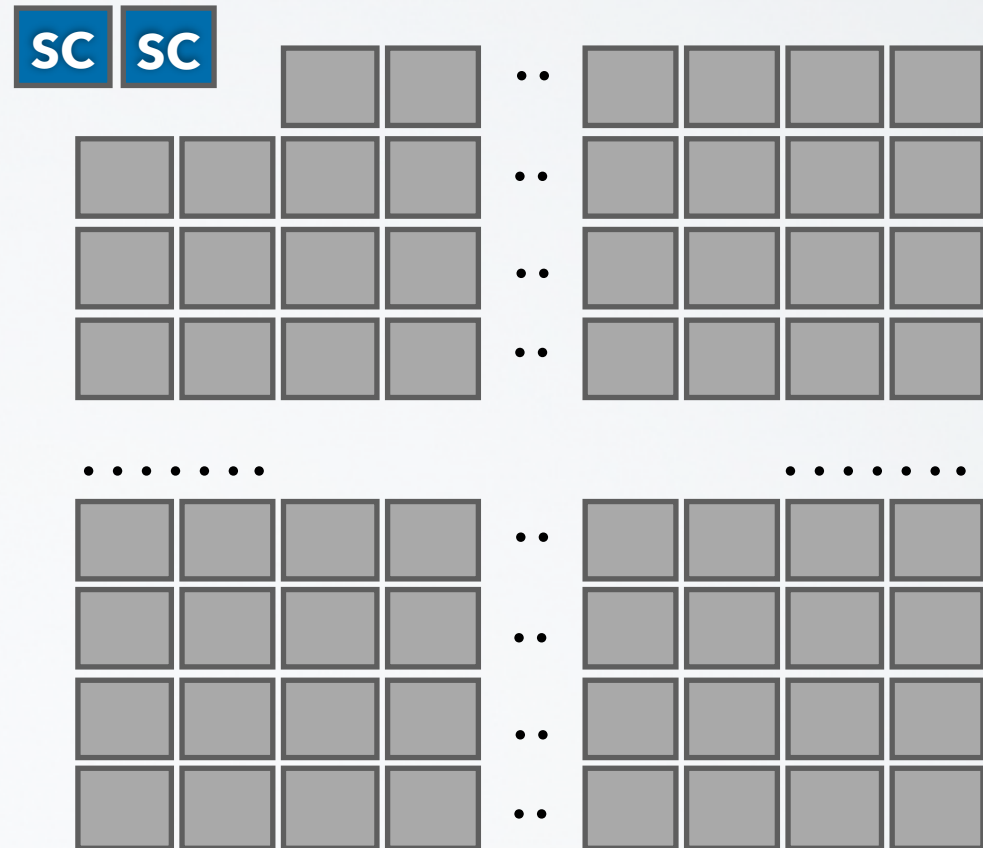


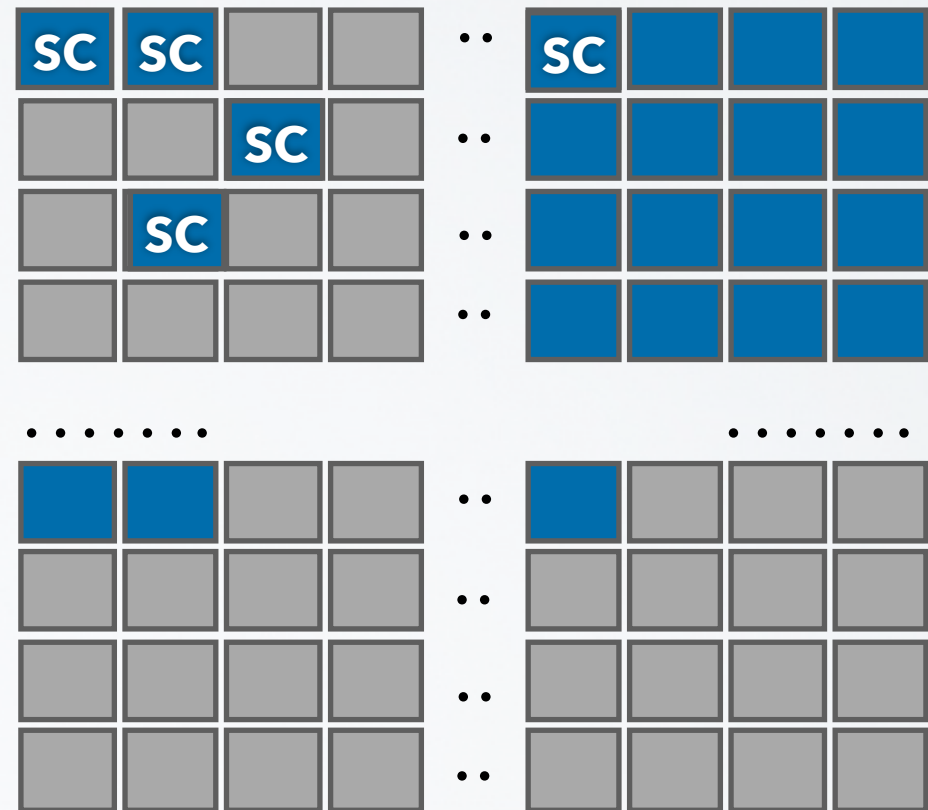




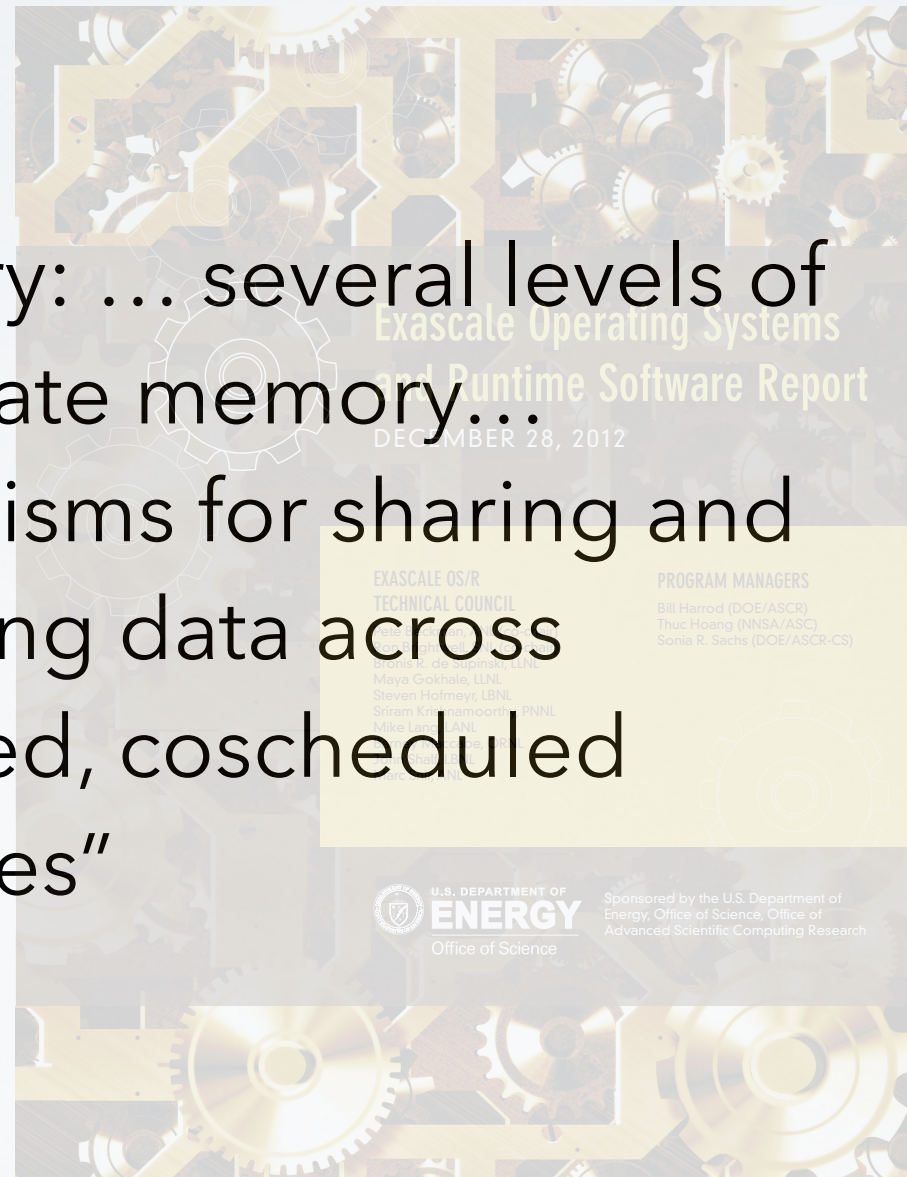


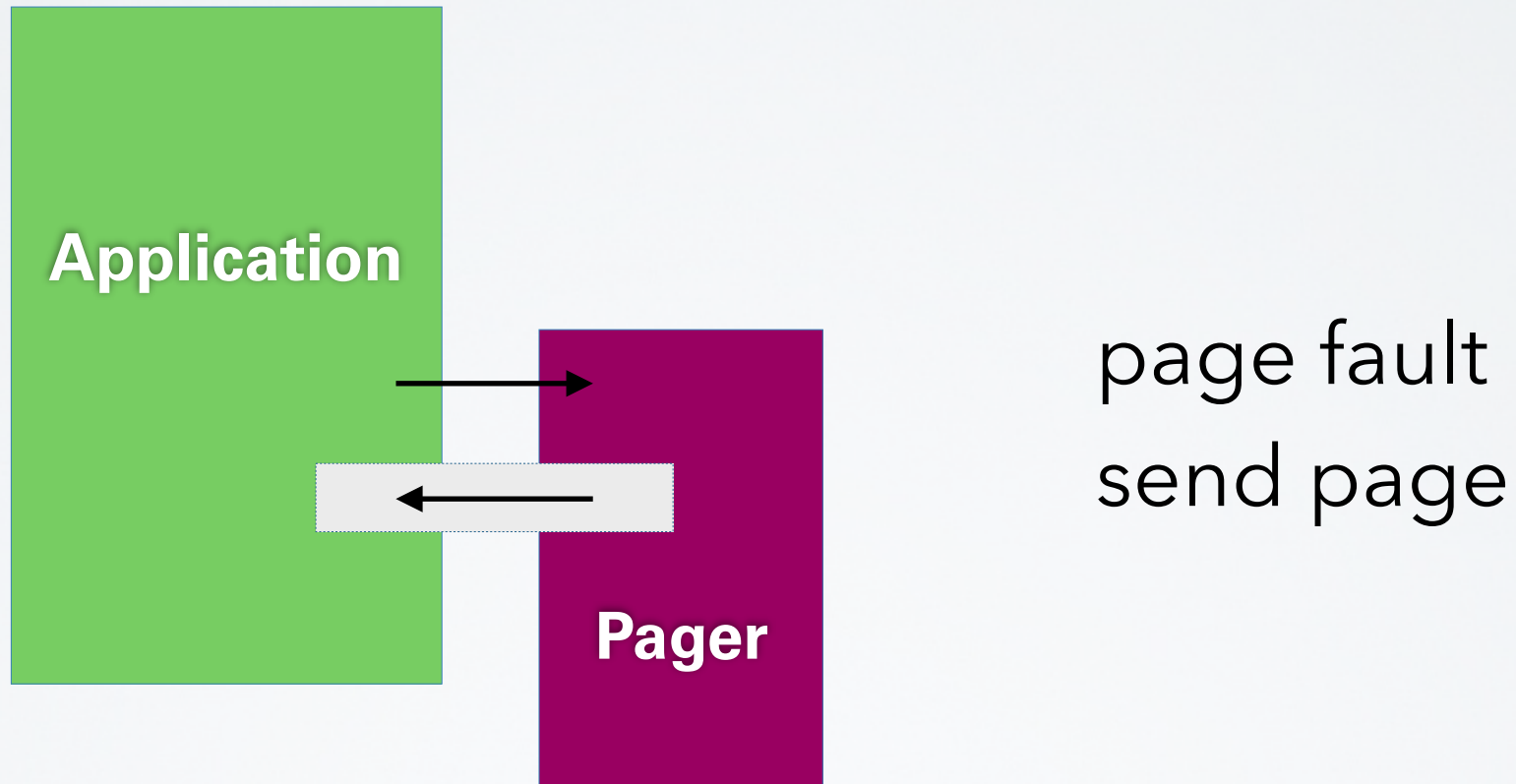
L4 Microkernel



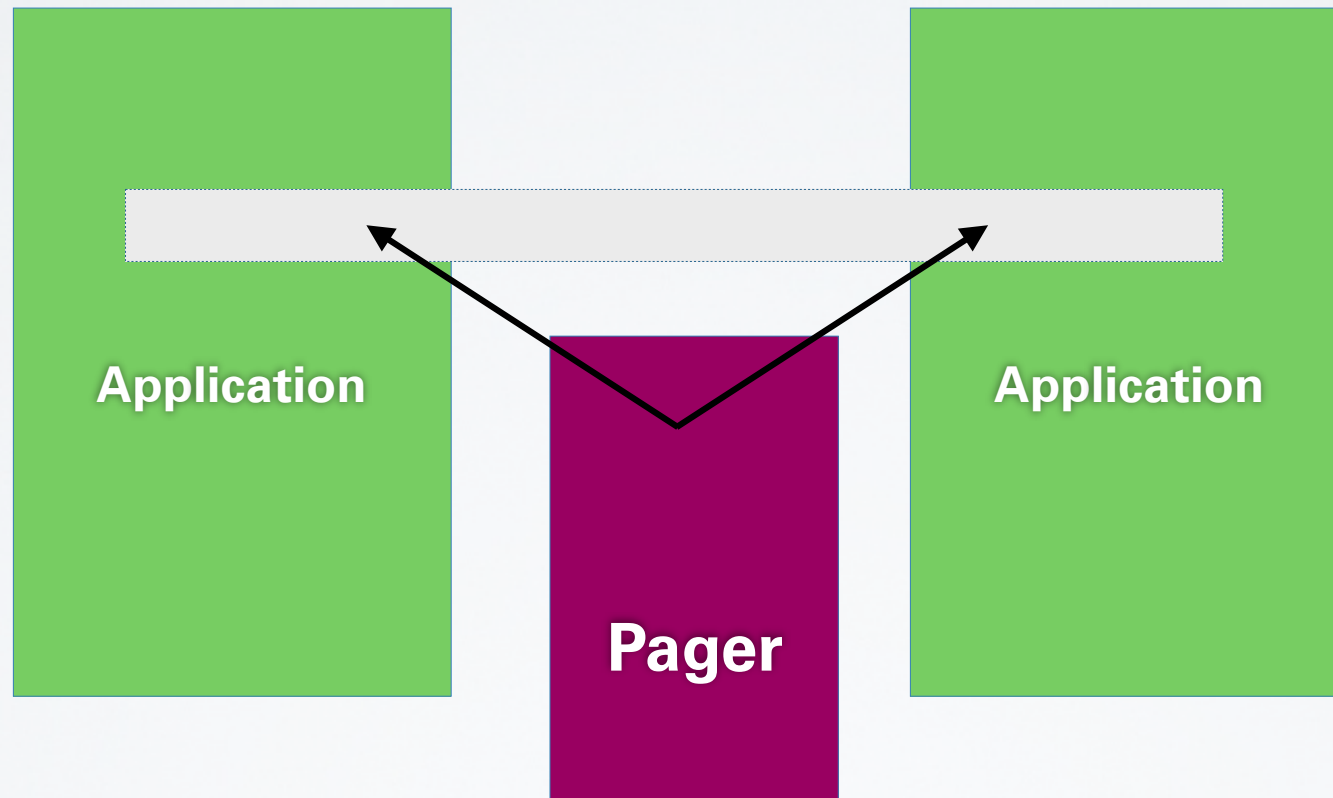


"Memory: ... several levels of solid-state memory... mechanisms for sharing and protecting data across colocated, coscheduled processes"



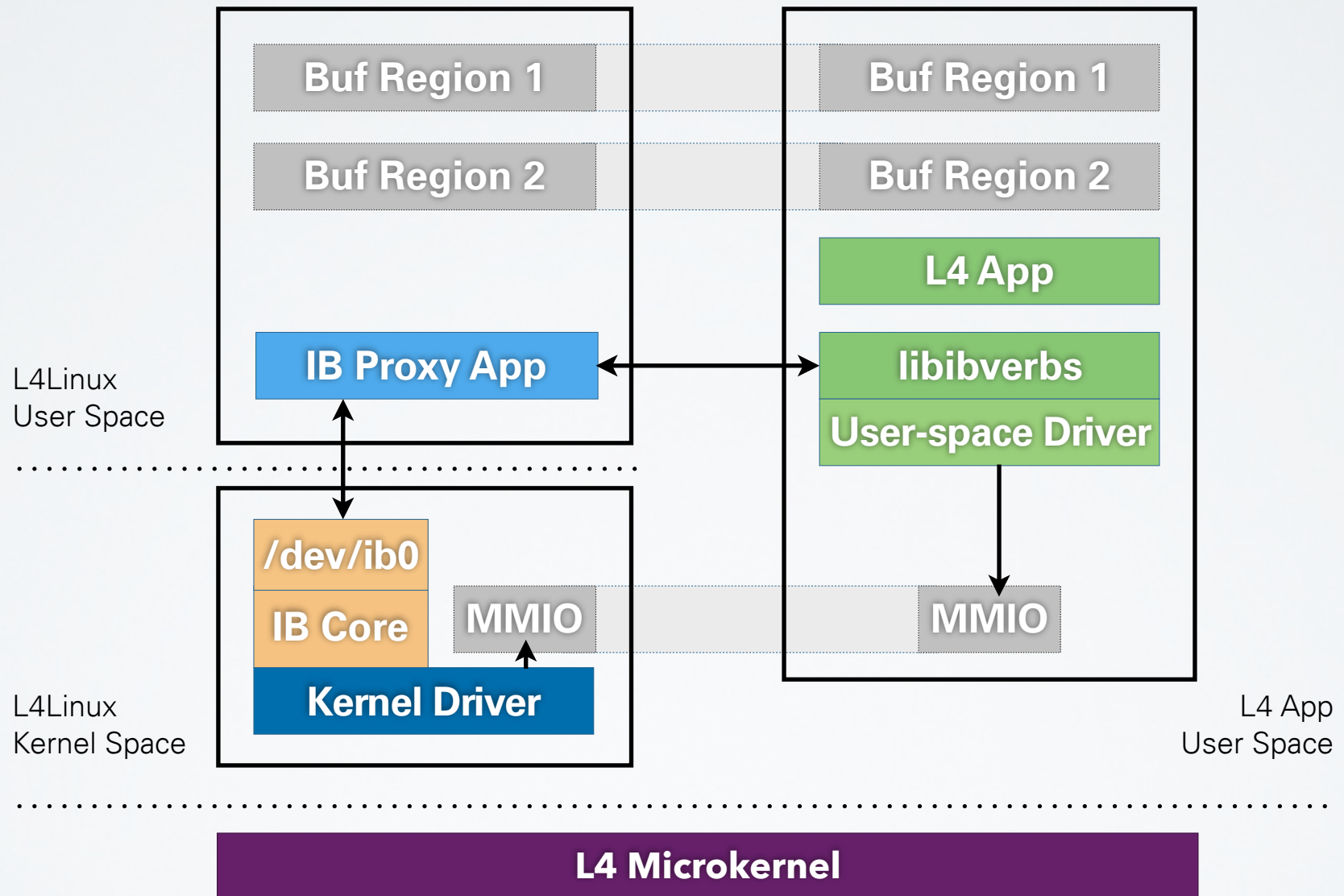


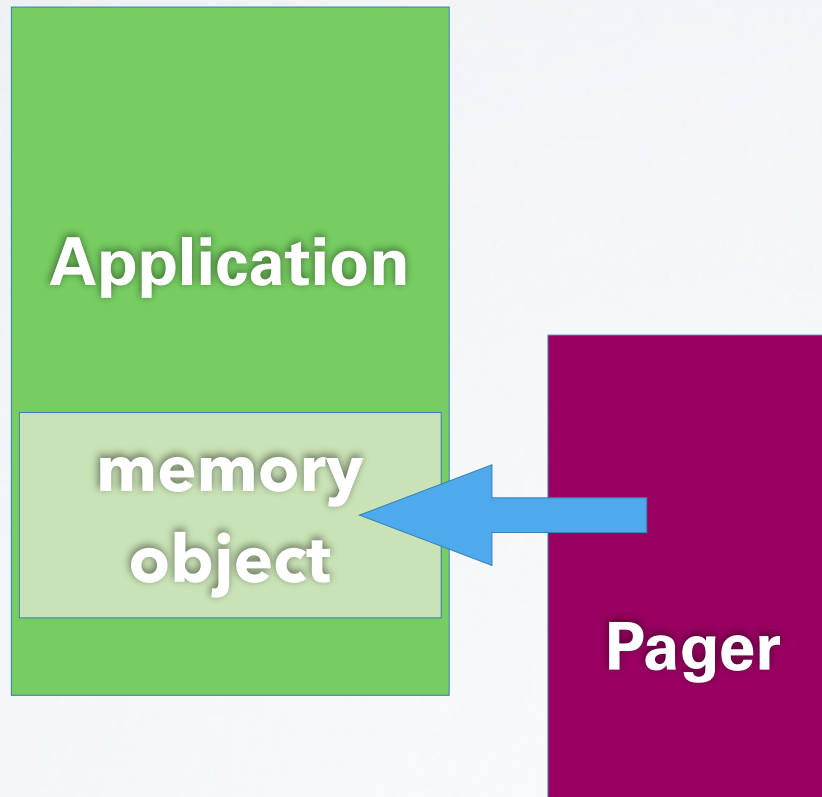
L4 Microkernel



L4 Microkernel

SPLIT INFINIBAND





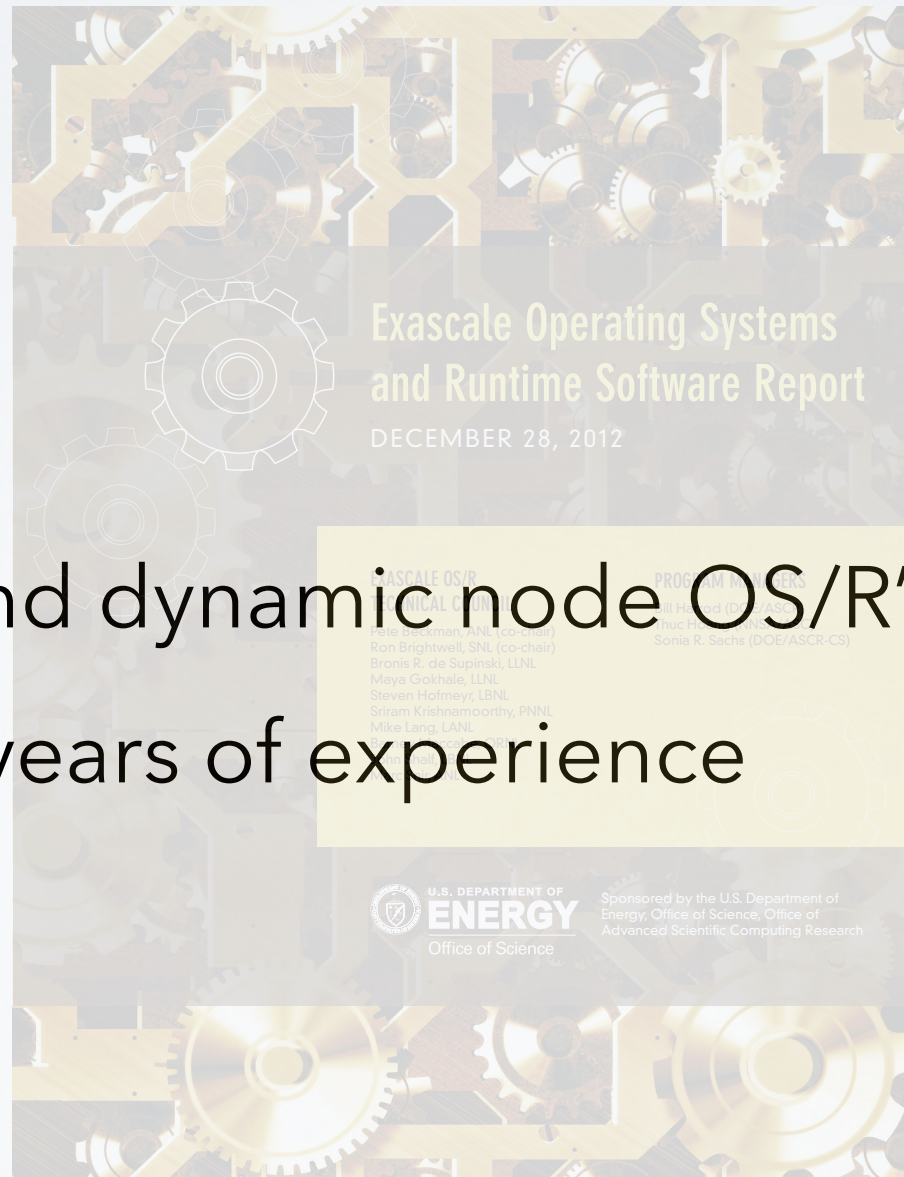
L4 Microkernel

interfacing with app and smart run-time

- division of work with smart run-time (e.g., Charm++ decomposition, balancing)
- impact on communication latency
- application level knowledge
resource usage, communication affinity

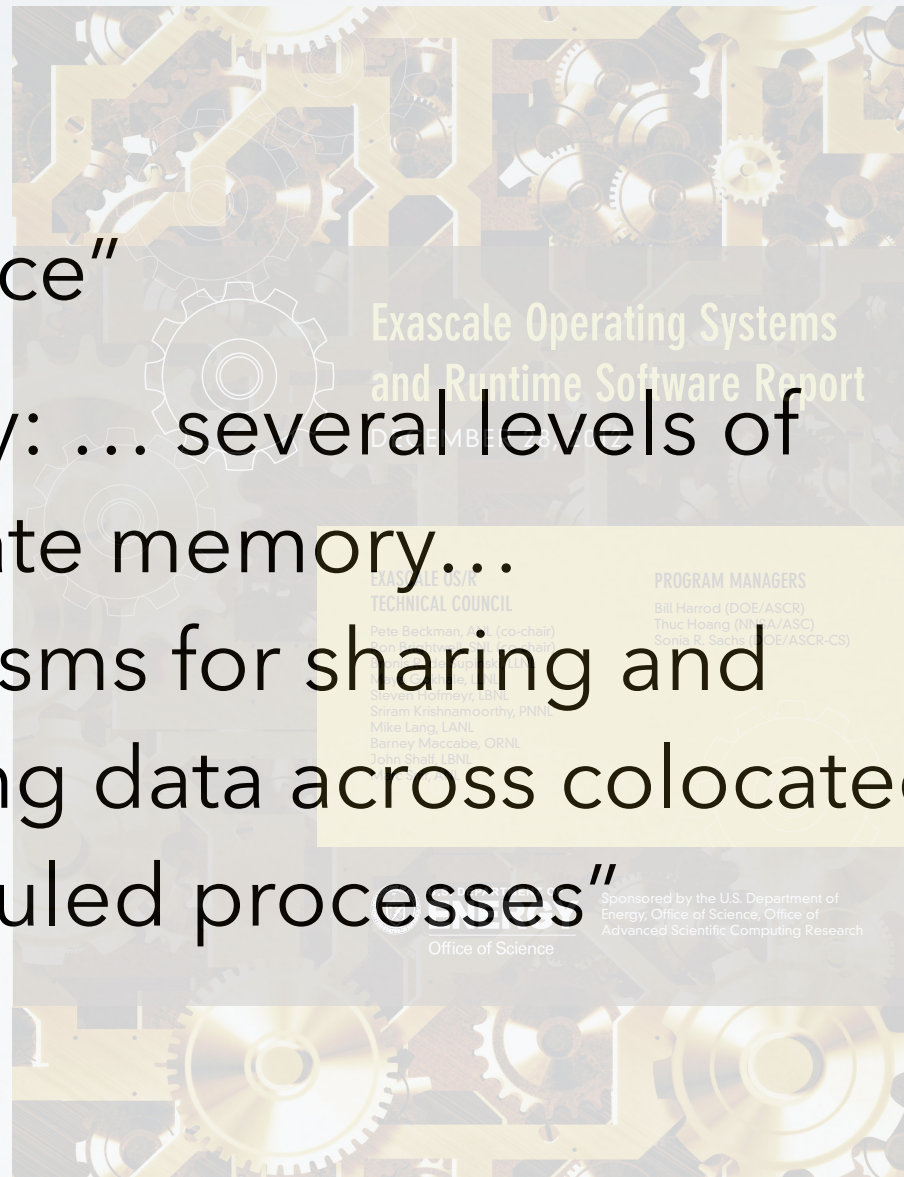
basis for

“agile and dynamic node OS/R”
with 20 years of experience

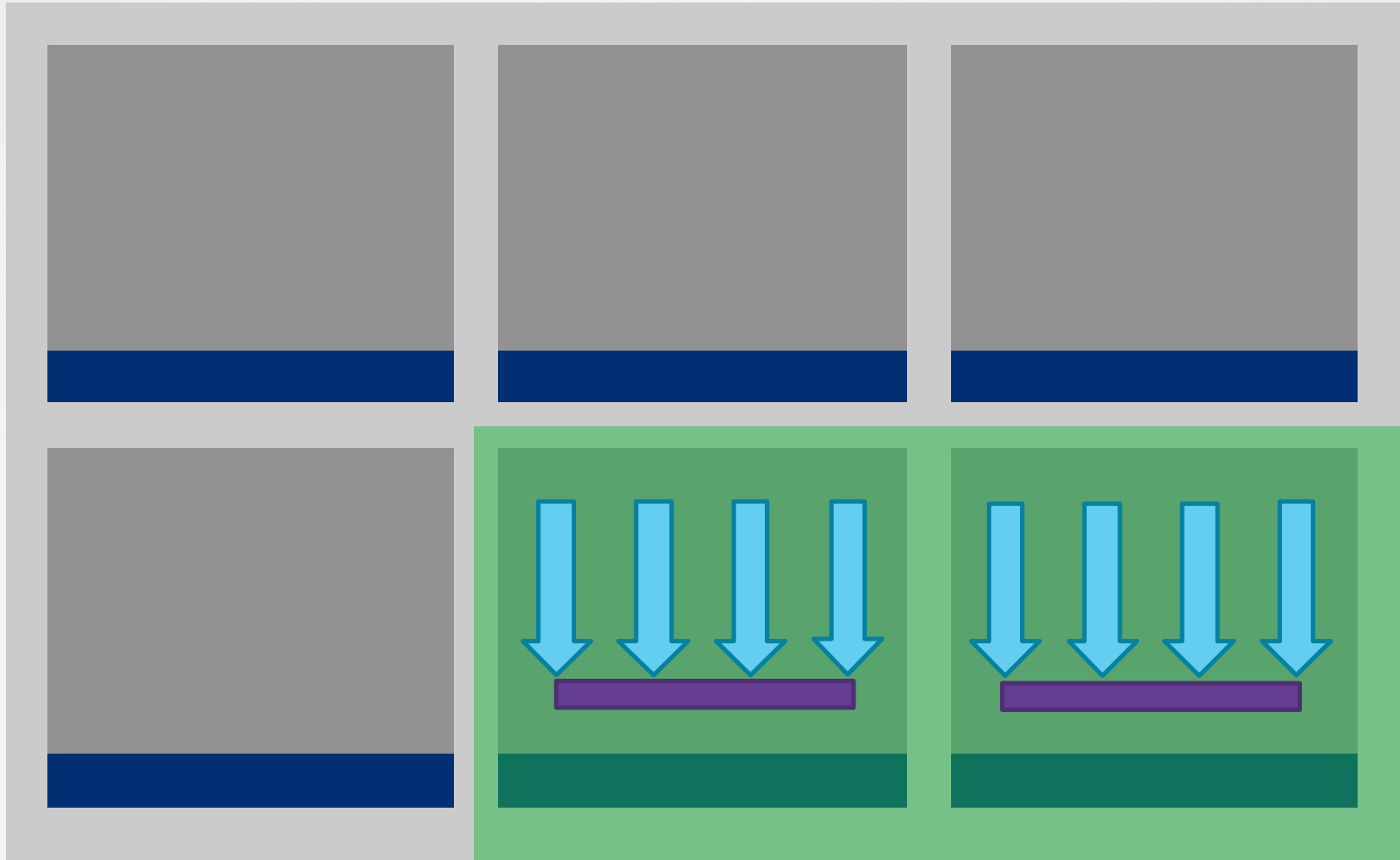


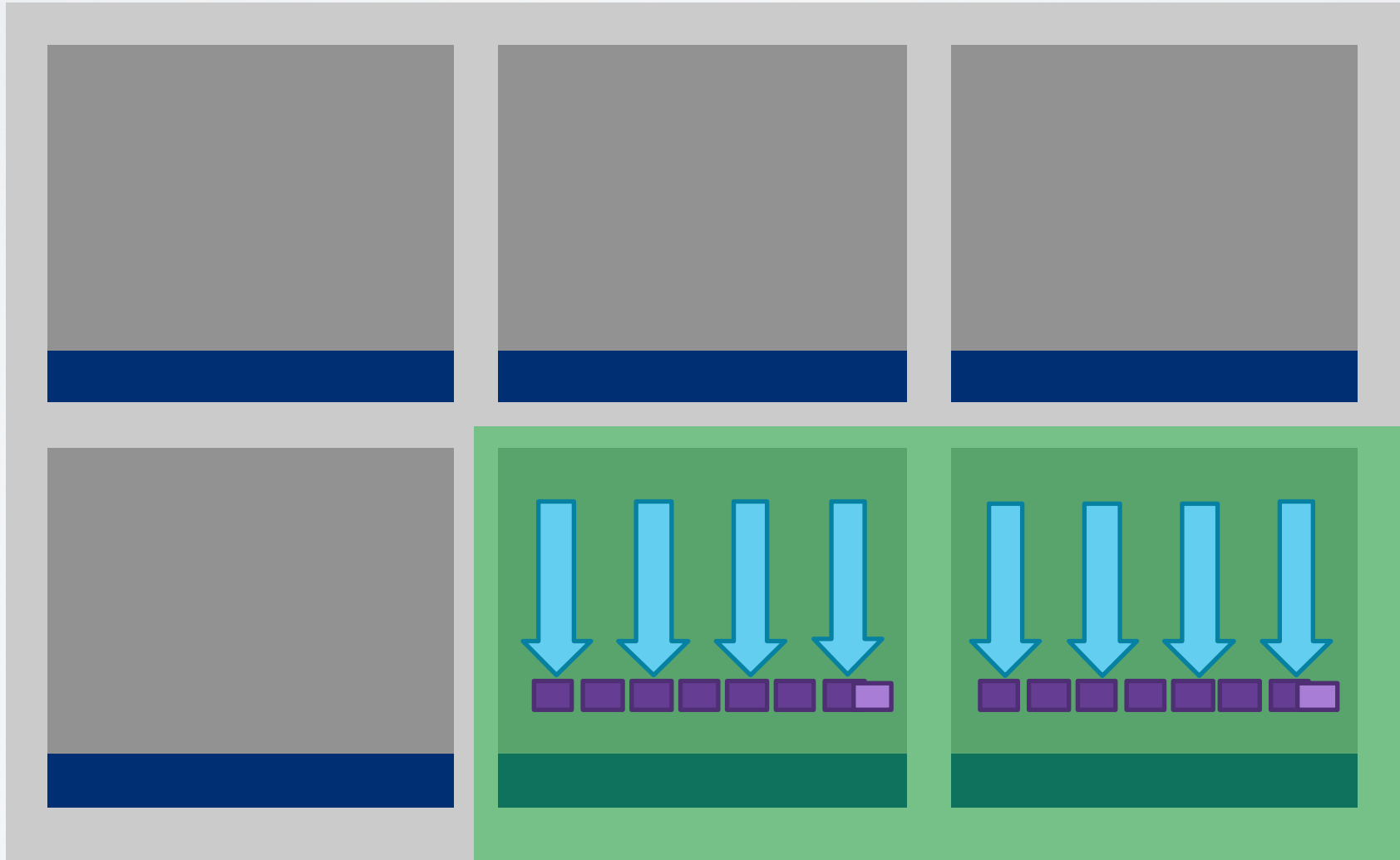
“Resilience”

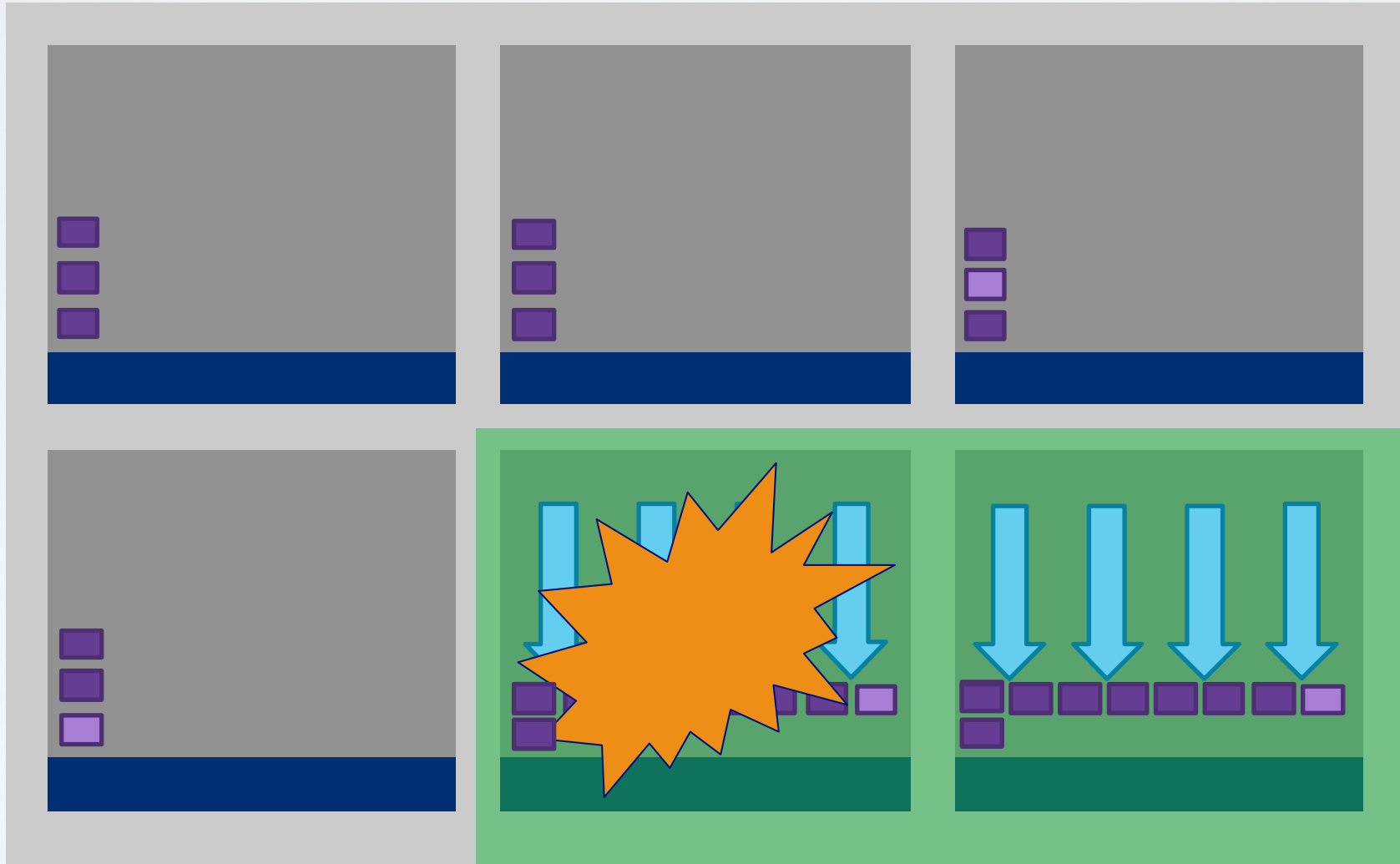
“Memory: ... several levels of solid- state memory... mechanisms for sharing and protecting data across colocated, coscheduled processes”

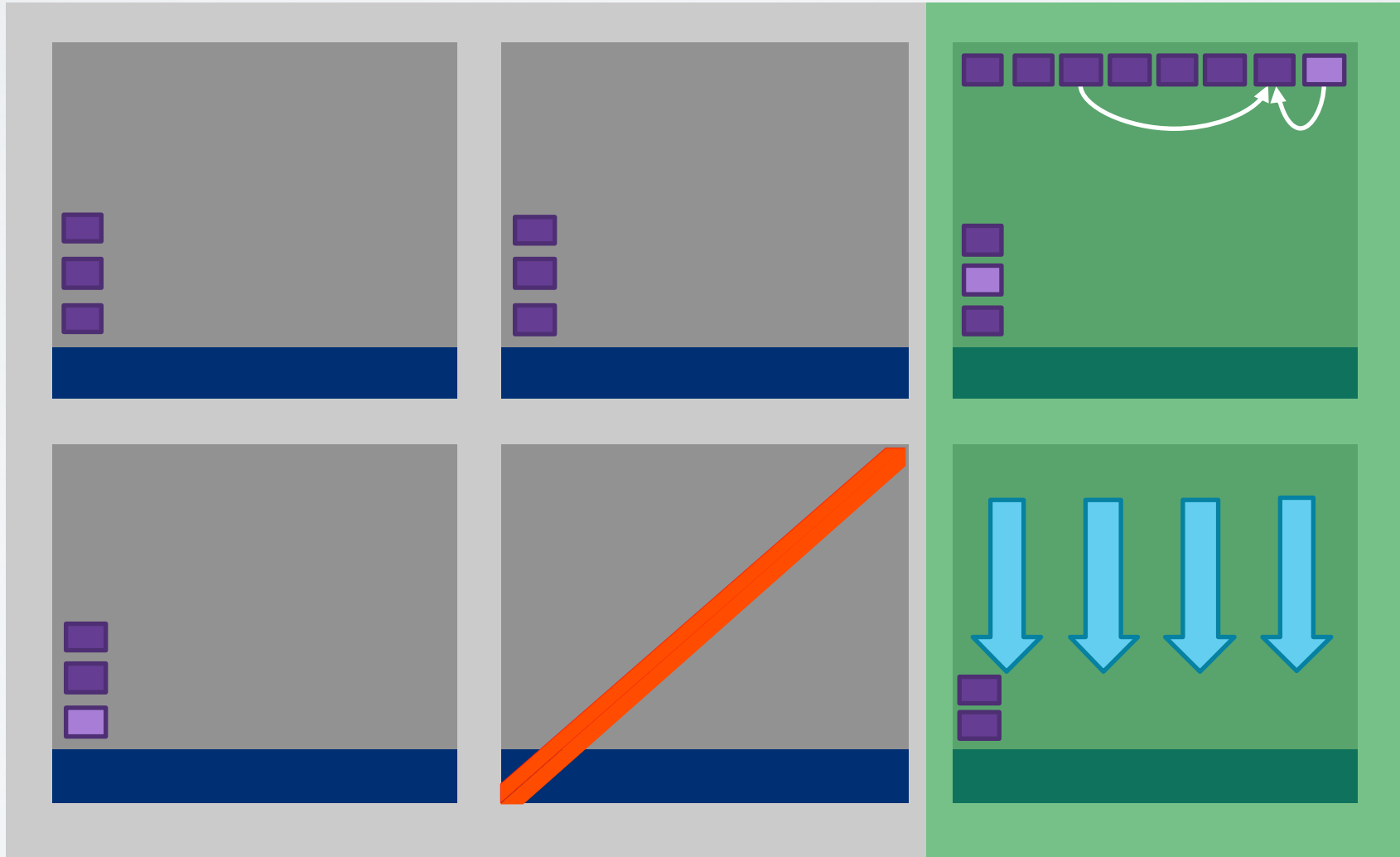


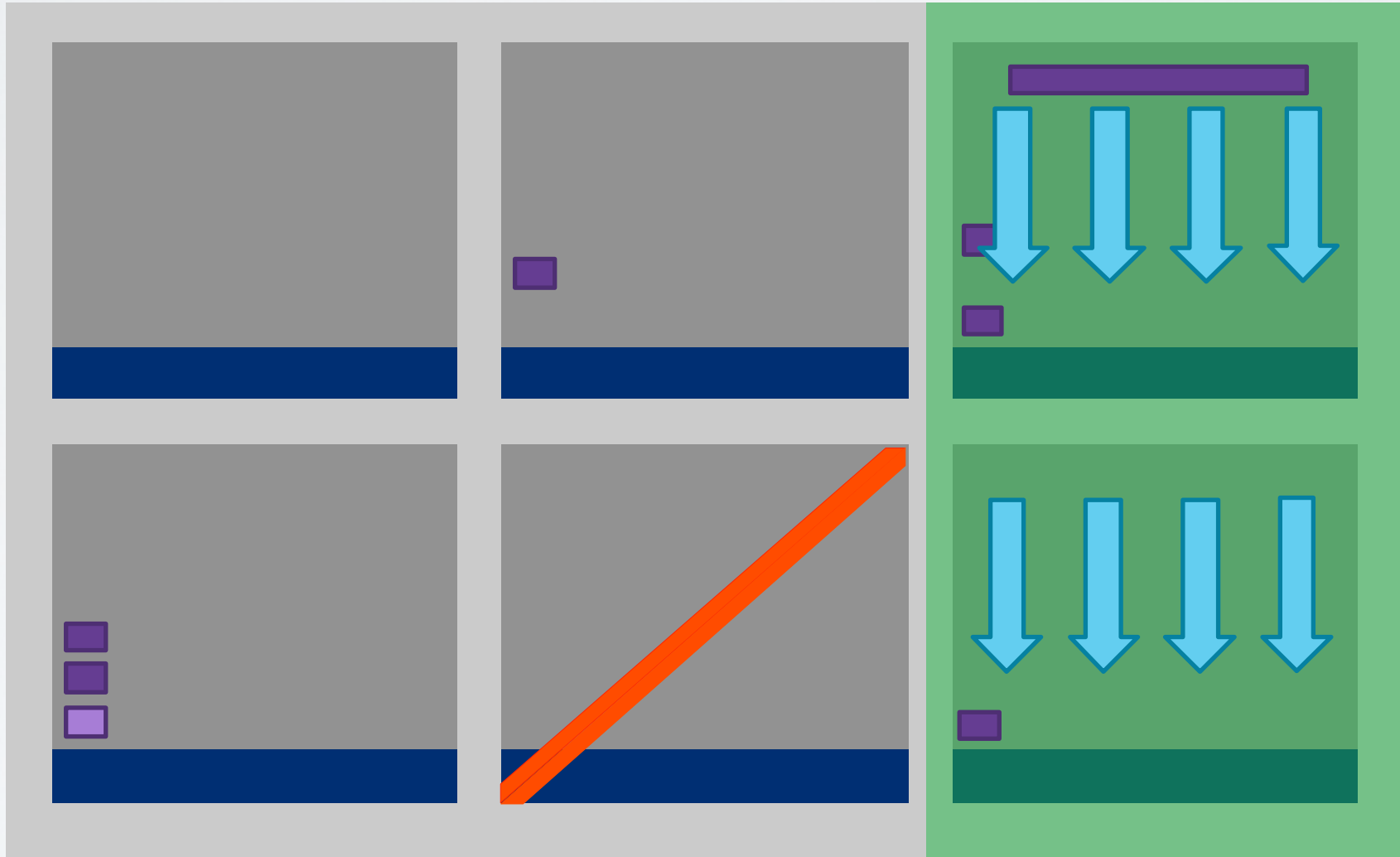
- higher failure rates
- diverse node-storage types
 - persistence
 - low-power DRAM

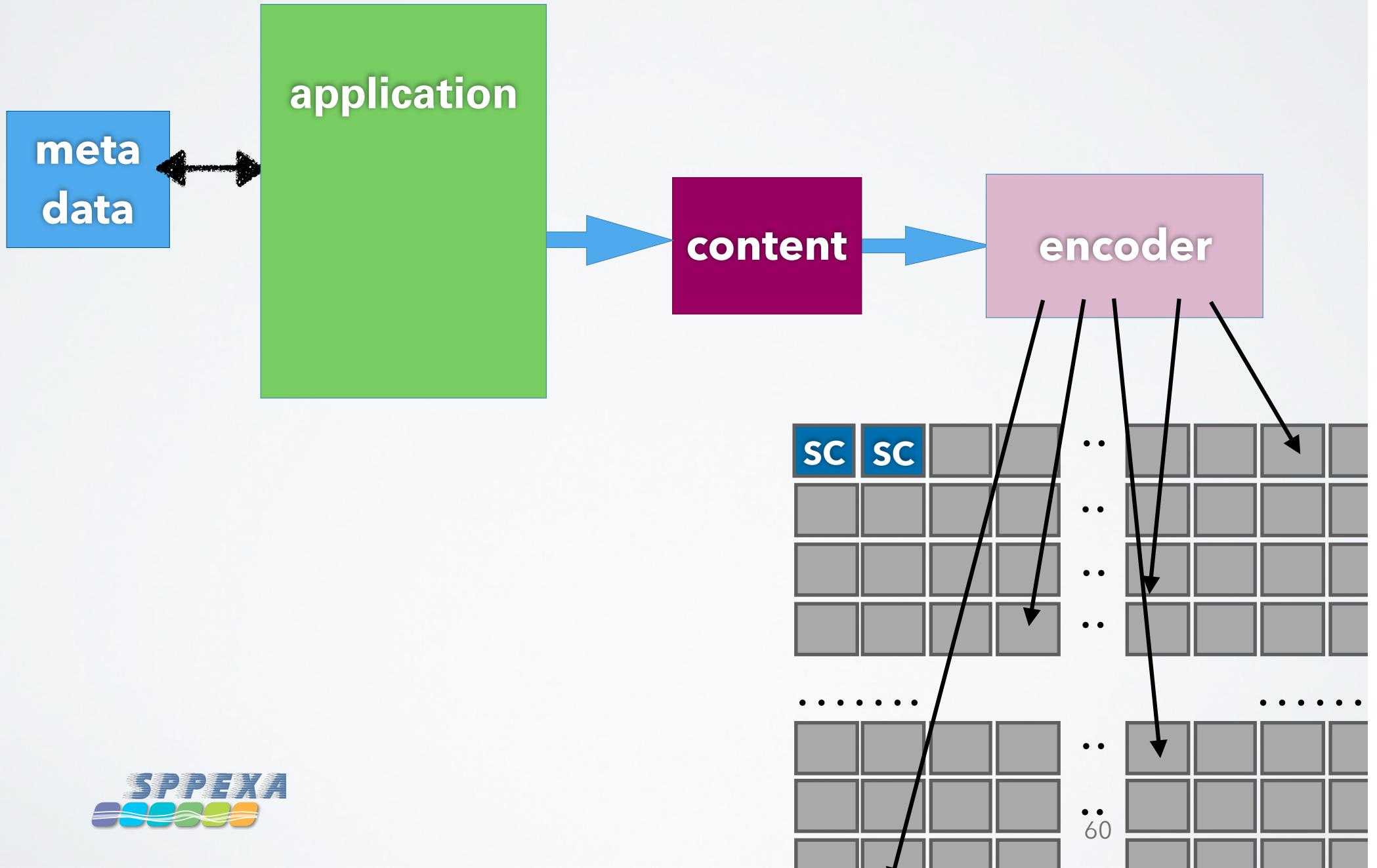


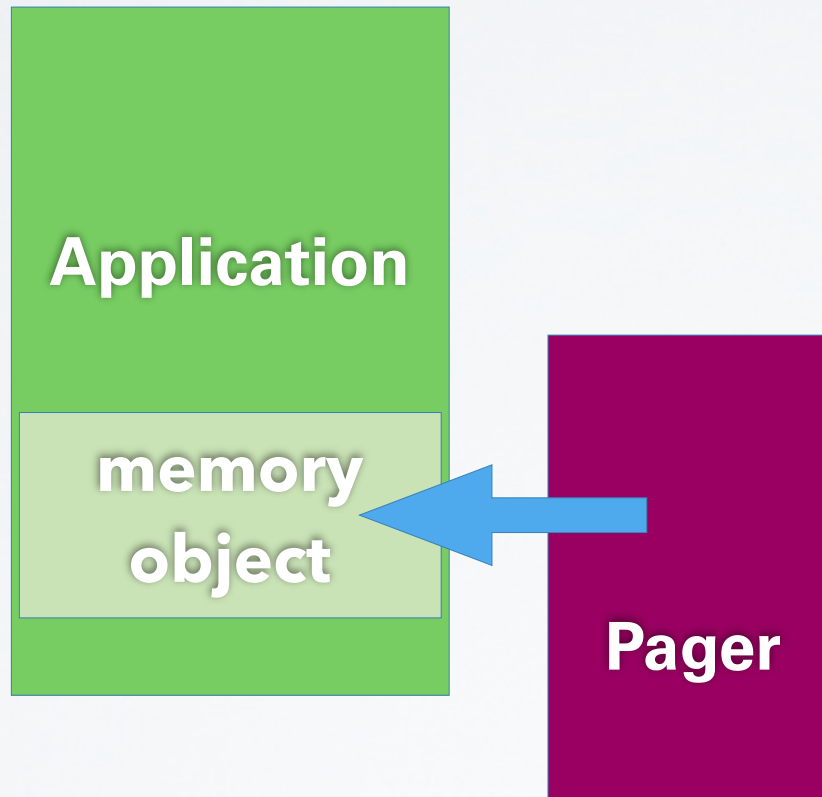




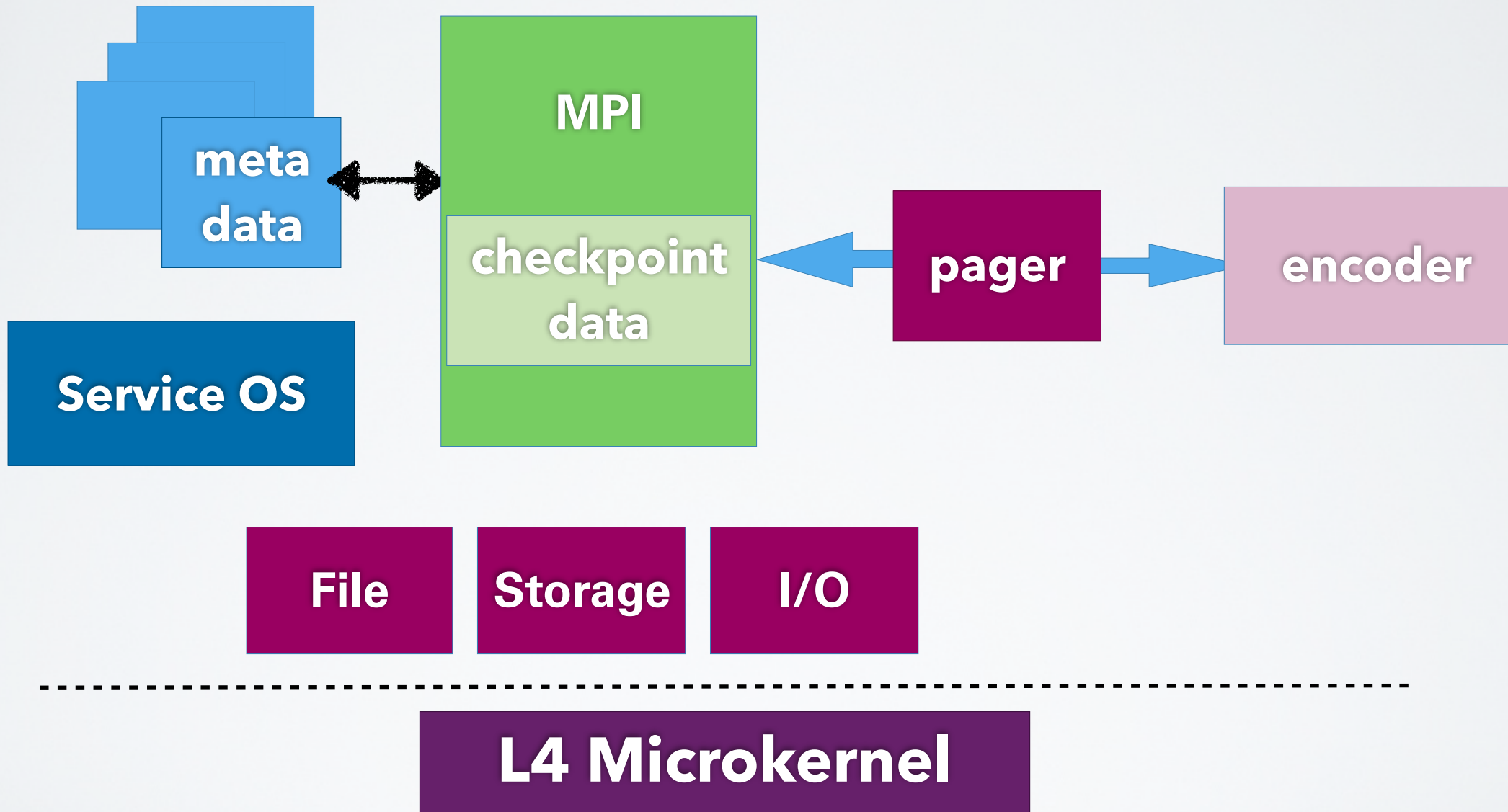






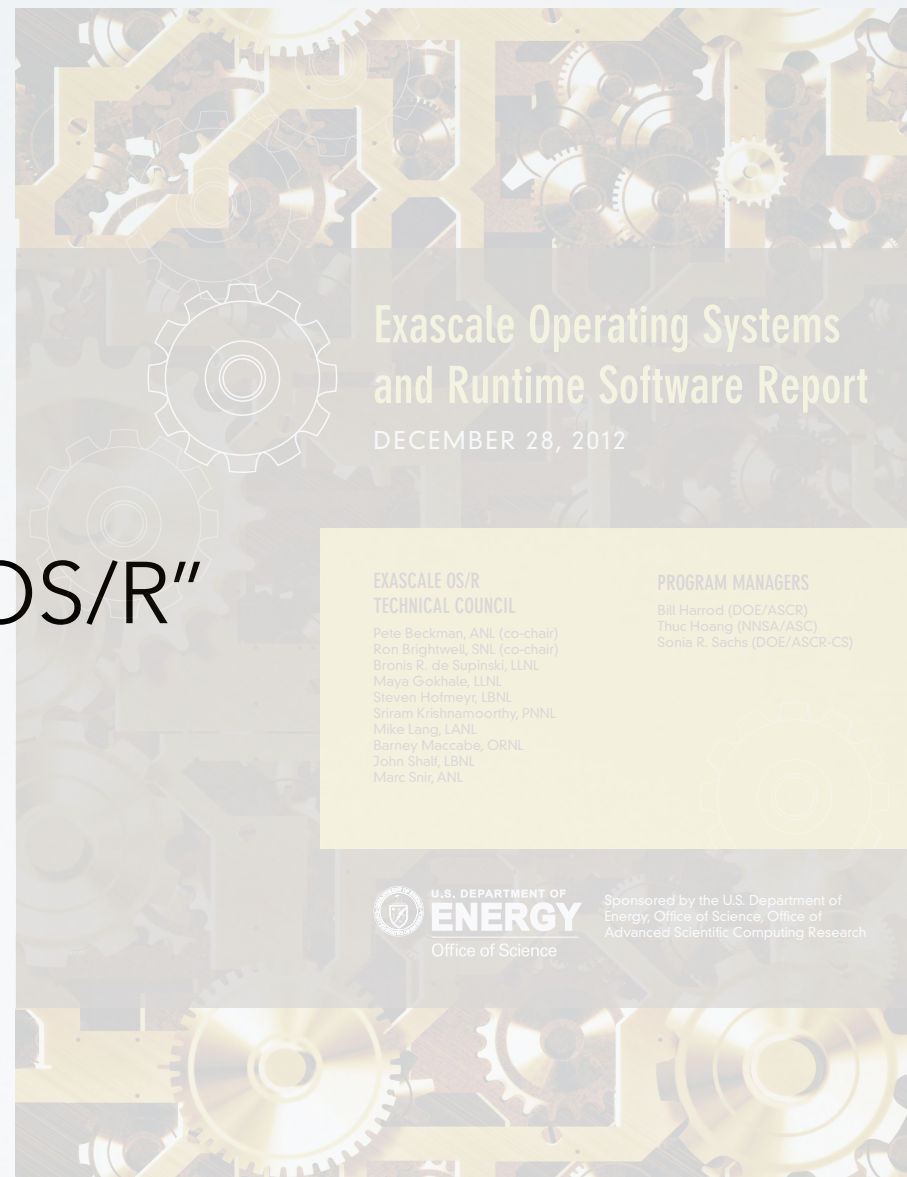


L4 Microkernel



- application ./ . run time ./ . OS interface
- allocation of encoded chunks
- global information on failures

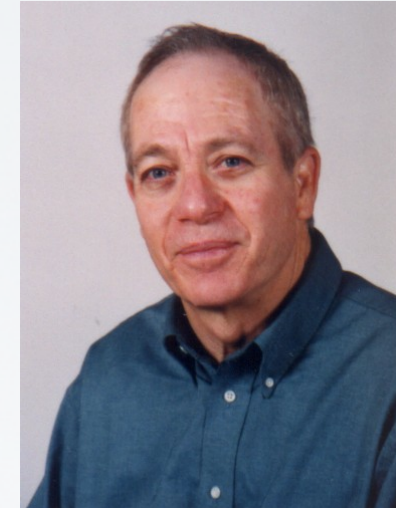
- “Global OS/R”



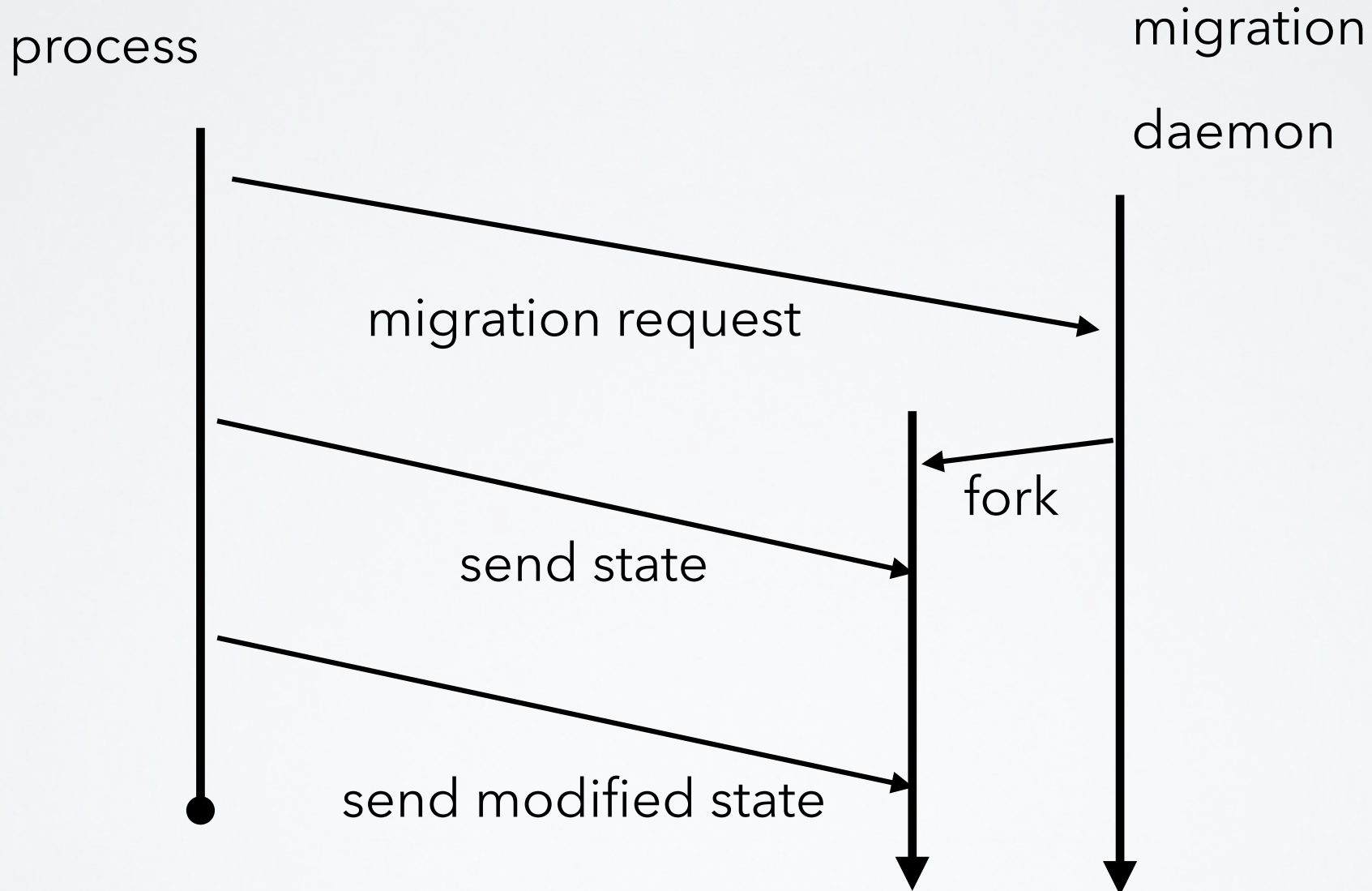
migration

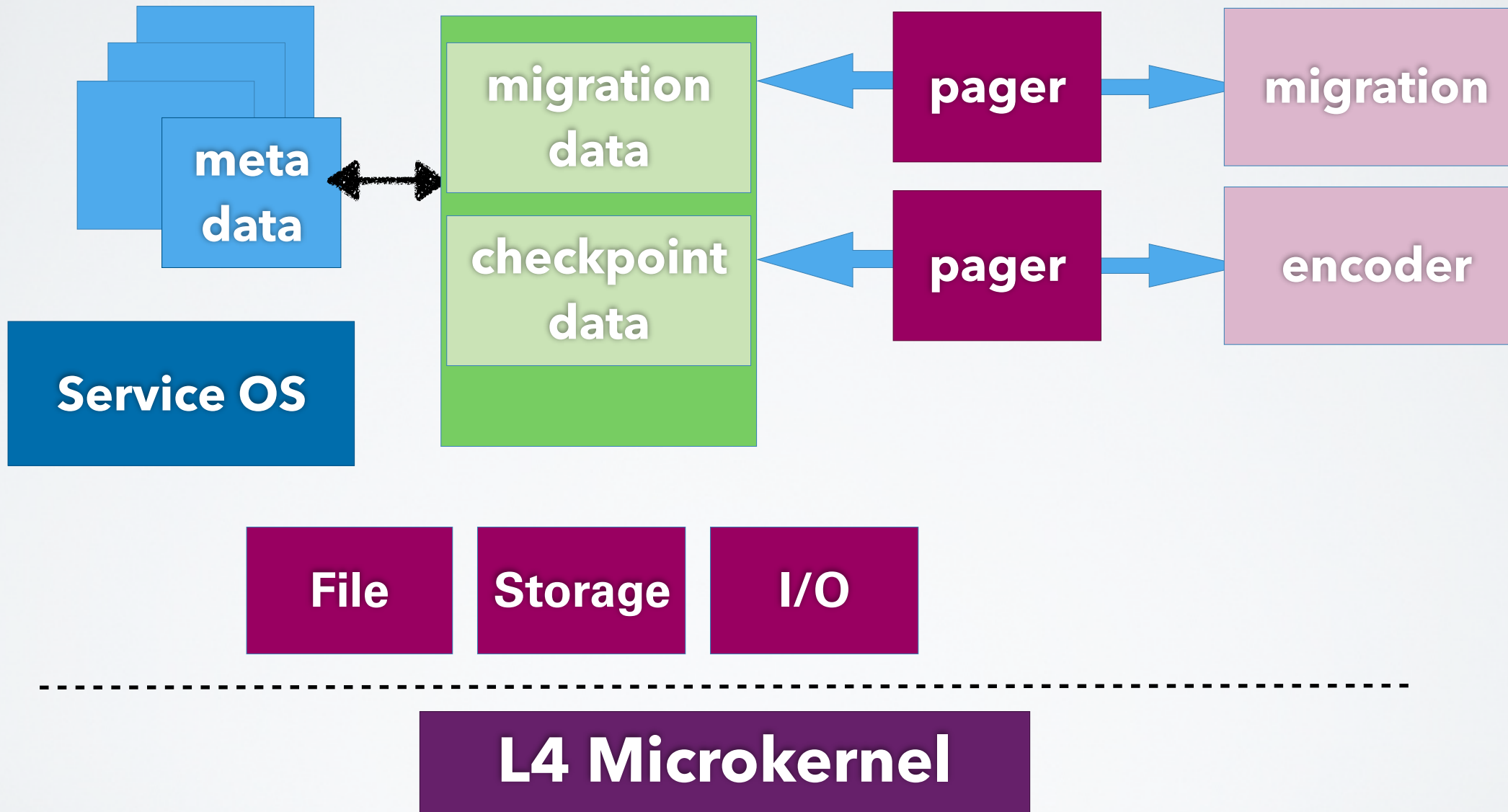
gossip

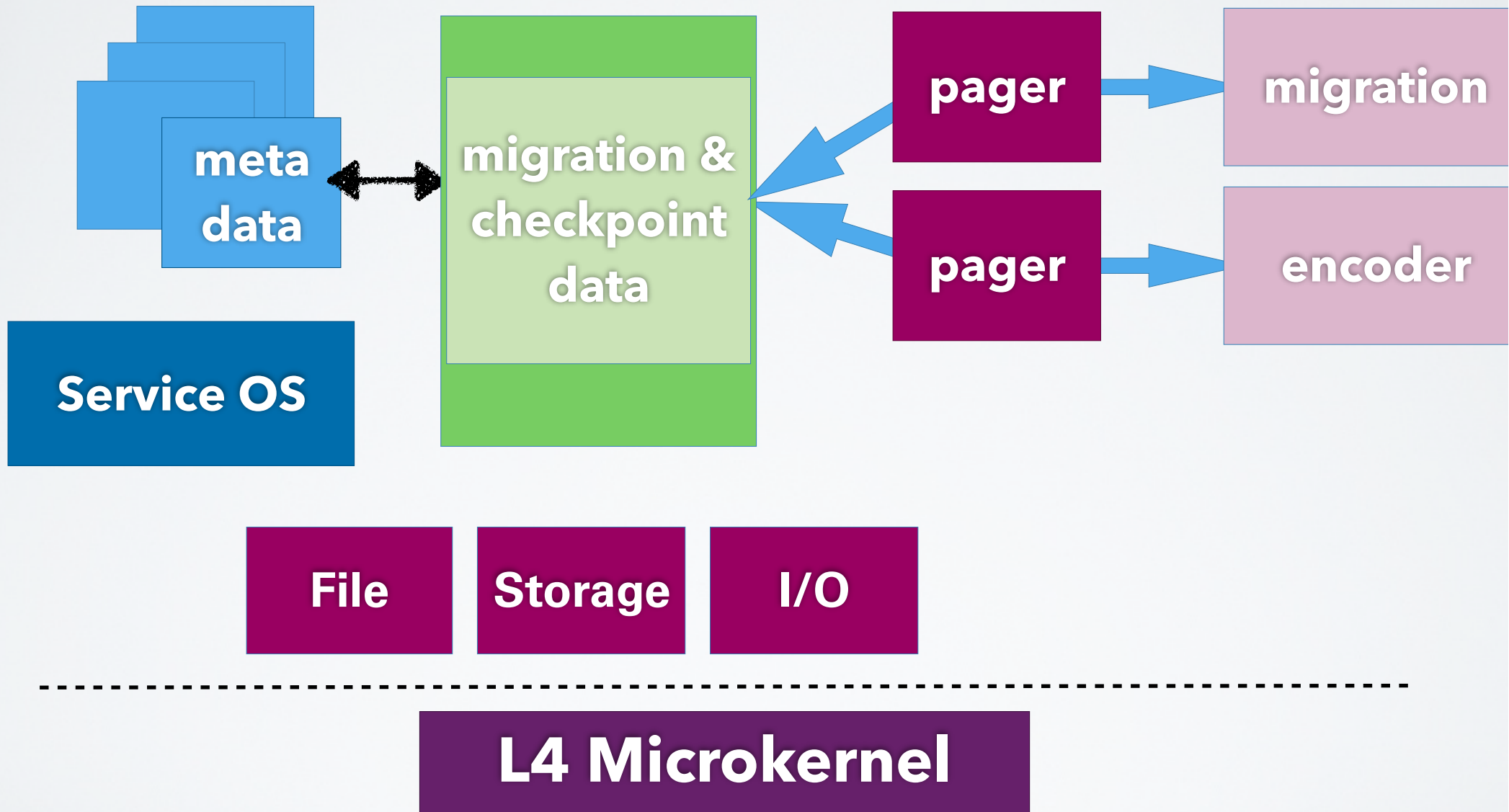
WWW

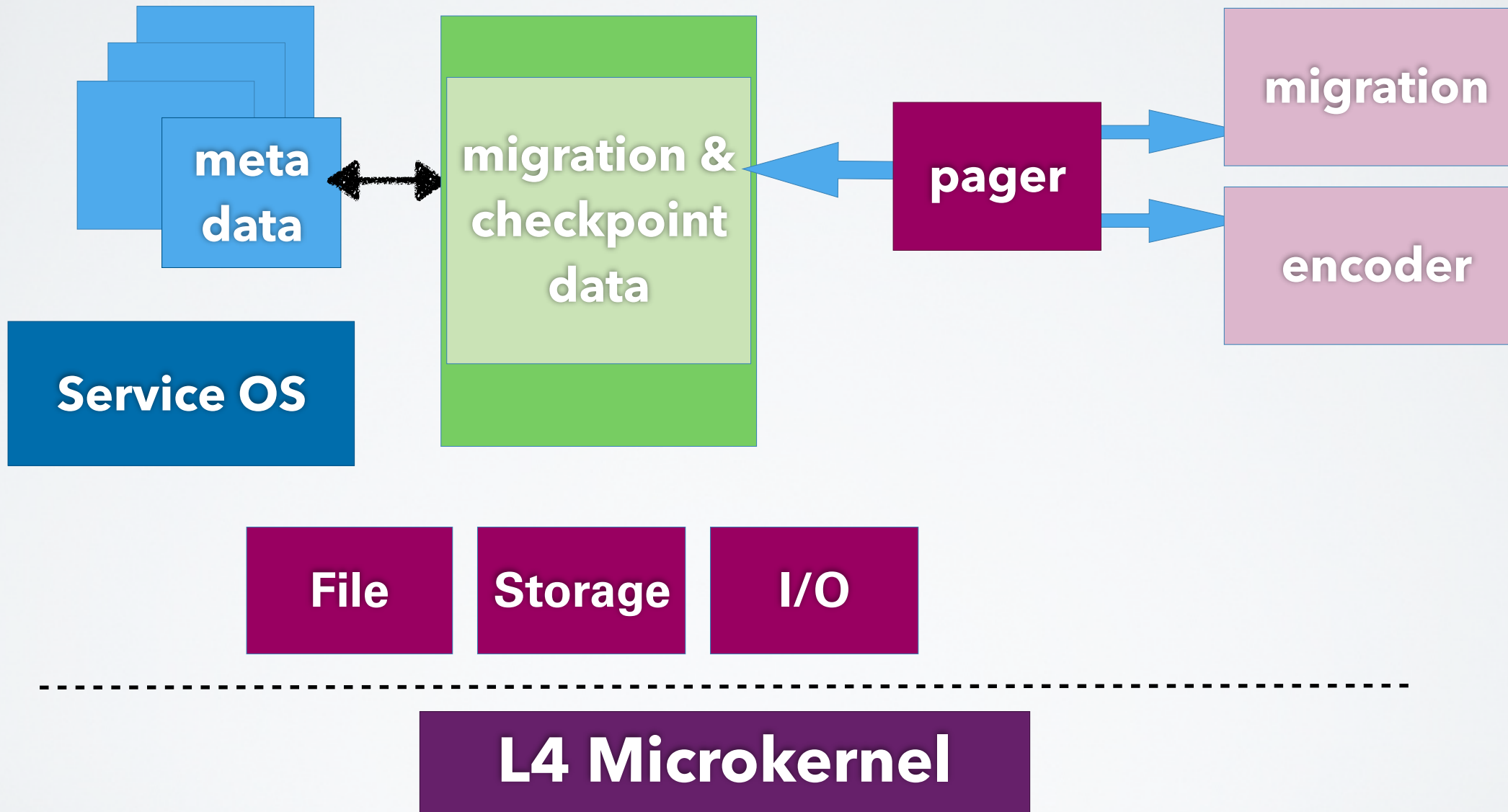


Amnon Barak

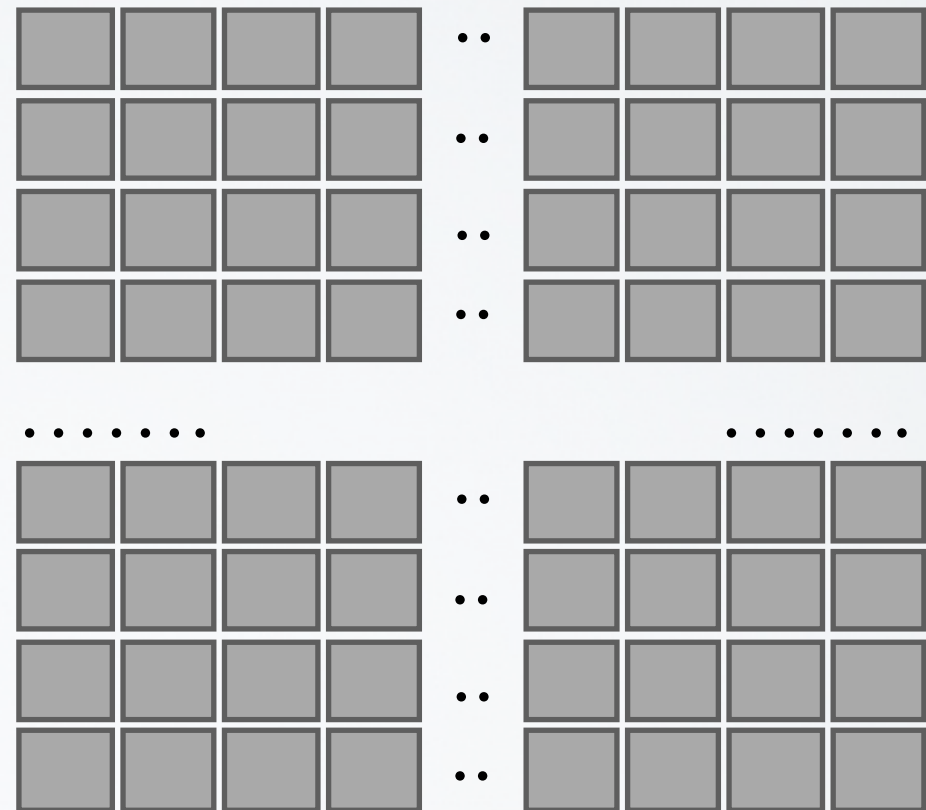
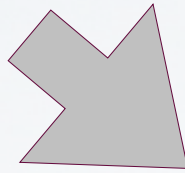








CENTRALIZED





Node 1

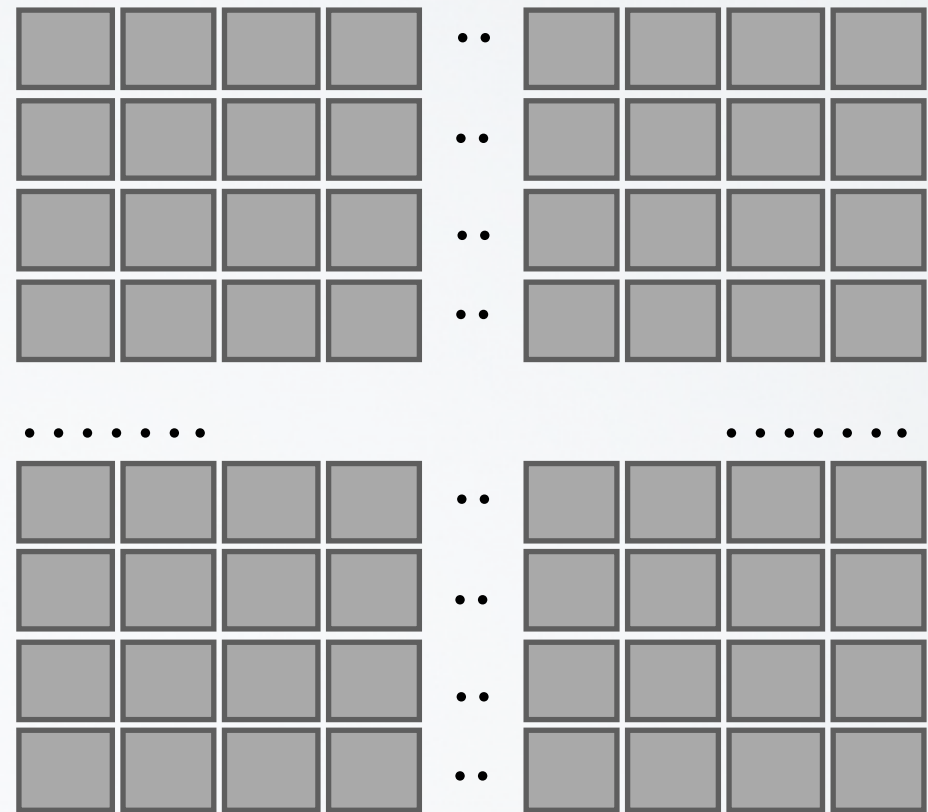
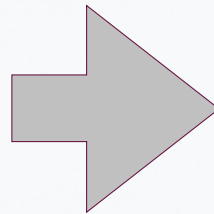
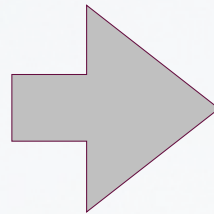
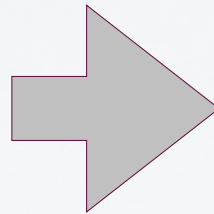


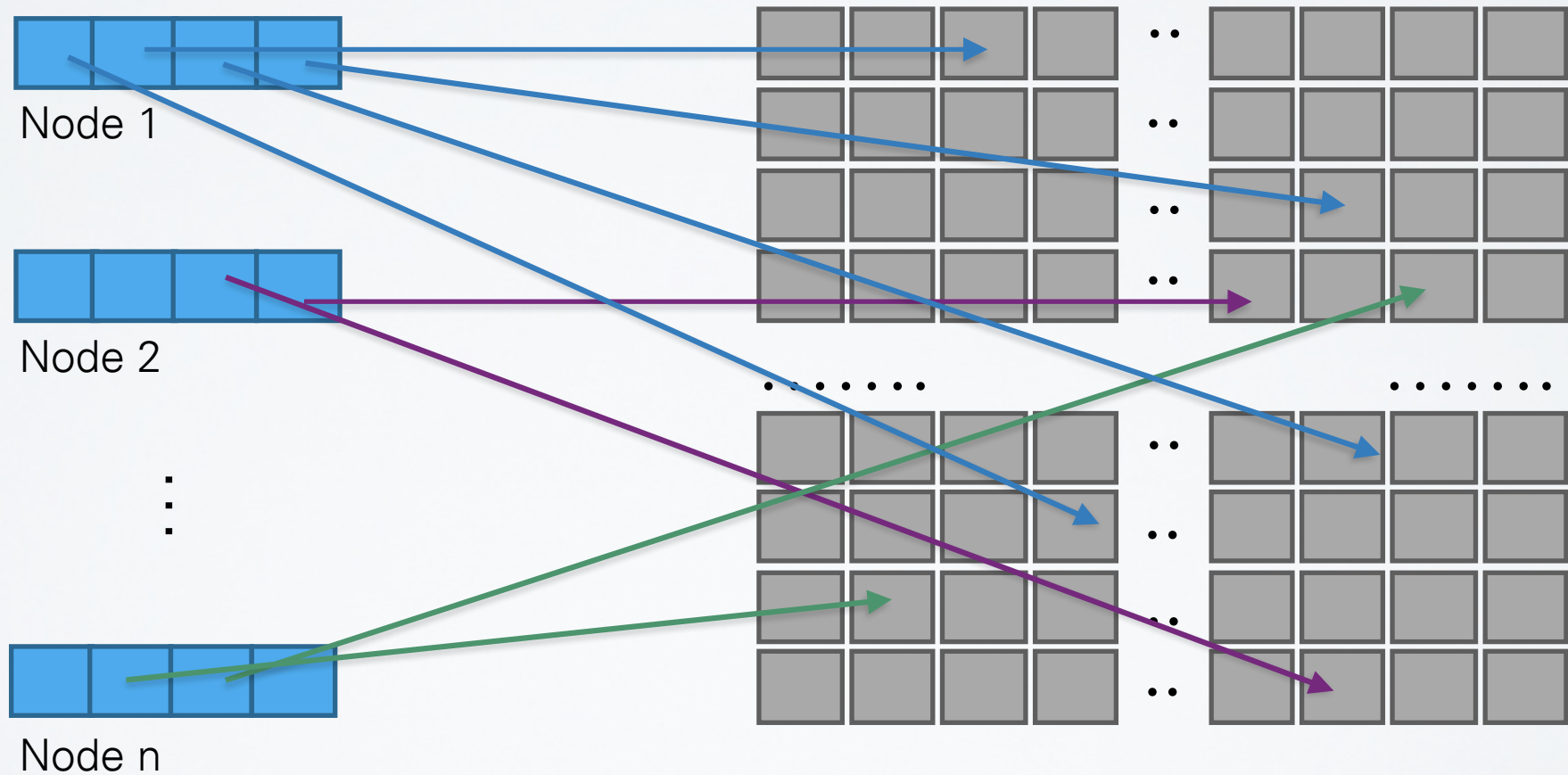
Node 2

⋮



Node n







Node 1

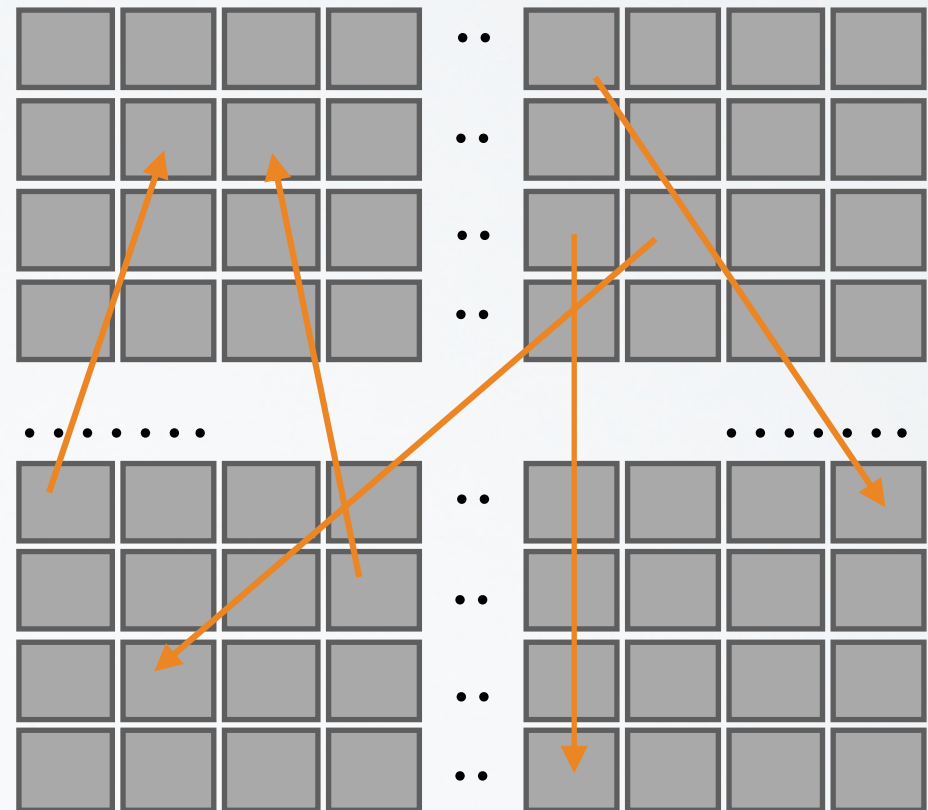


Node 2

⋮



Node n





Node 1



Node 2

⋮



Node n

When

M: load difference discovered
anomaly discovered
anticipated

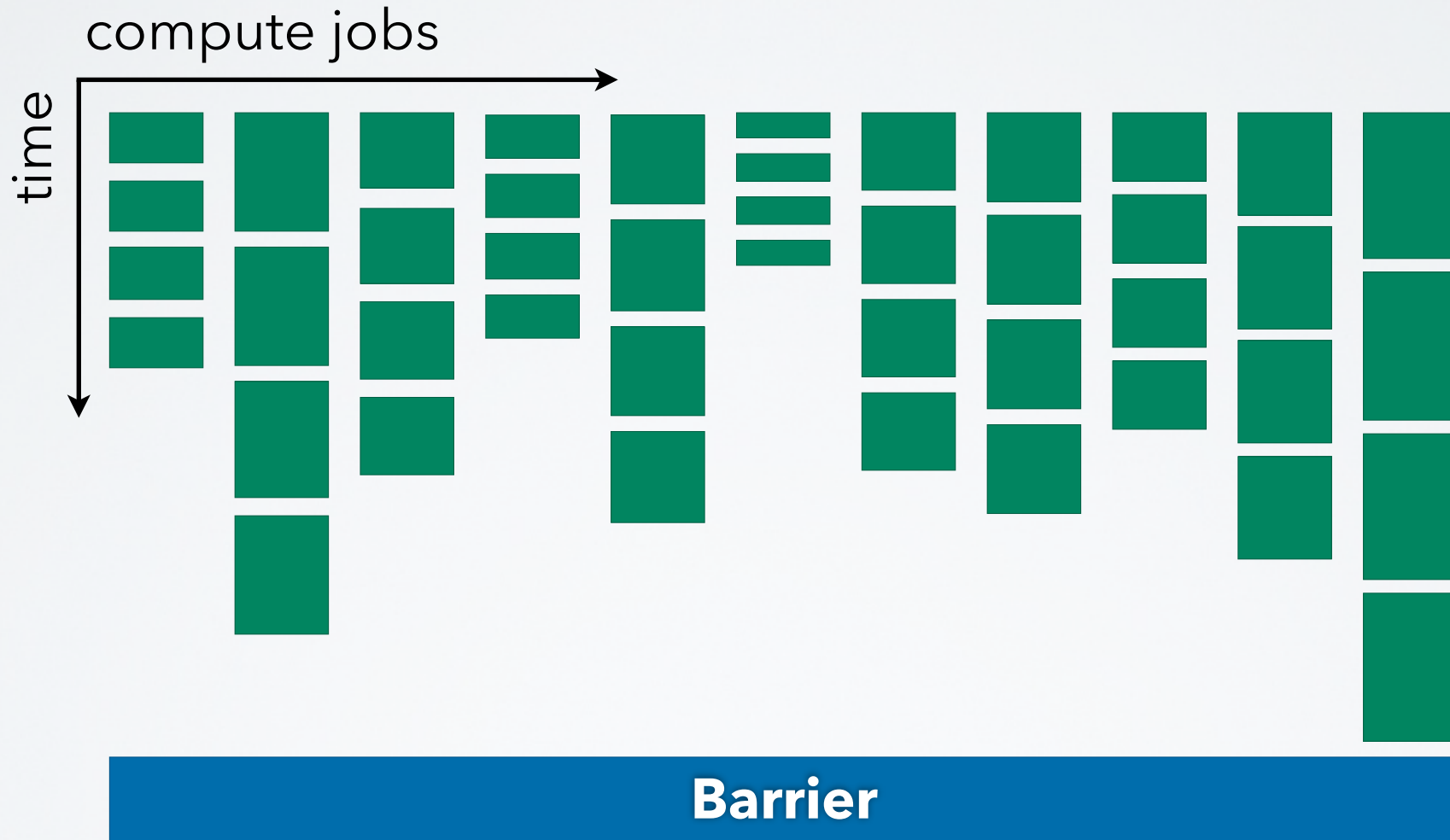
Where

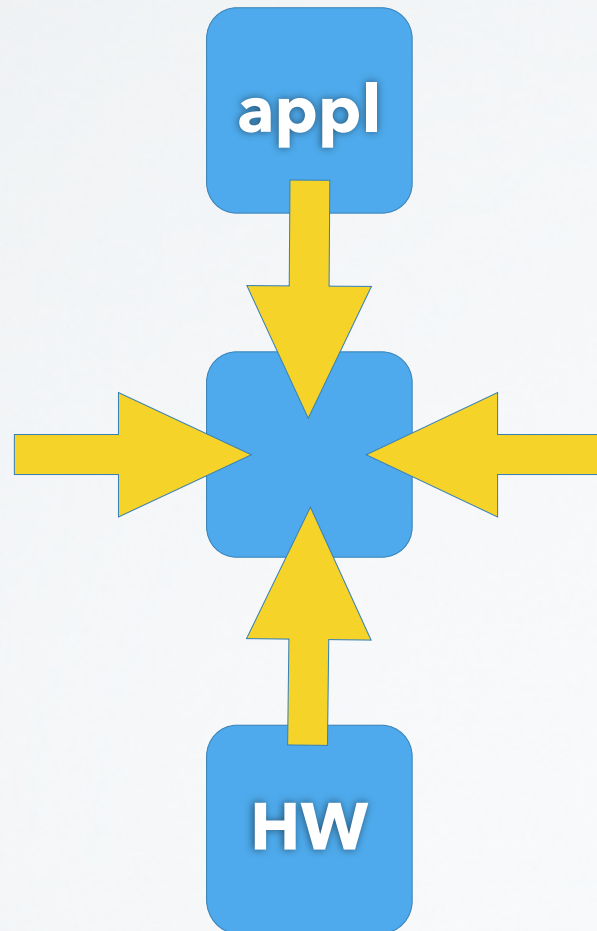
M: memory, cycles, comm
consider topology
application knowledge

Which

M: past predicts future
application knowledge

- gossip scalability
-> talk by Carsten Weinhold
- decentralized topology-aware allocation
- impact of migration on communication
- what is load ?
- how to “yell for help” ?



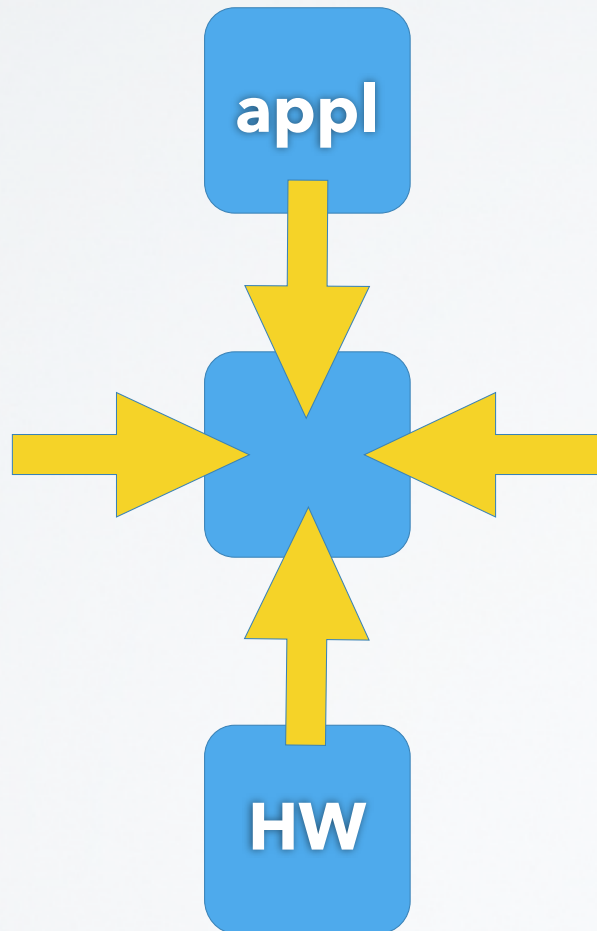


past predicts future

gossip

event counters

clocks



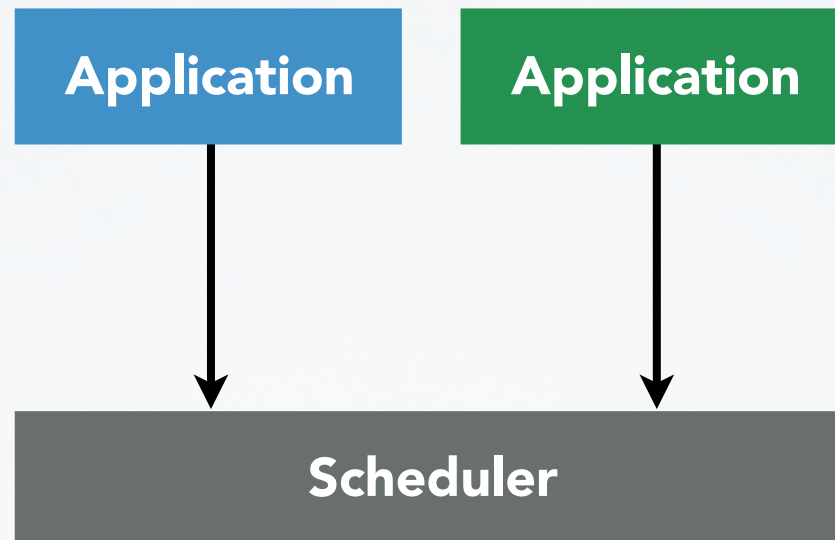
applications predict future

cloud density

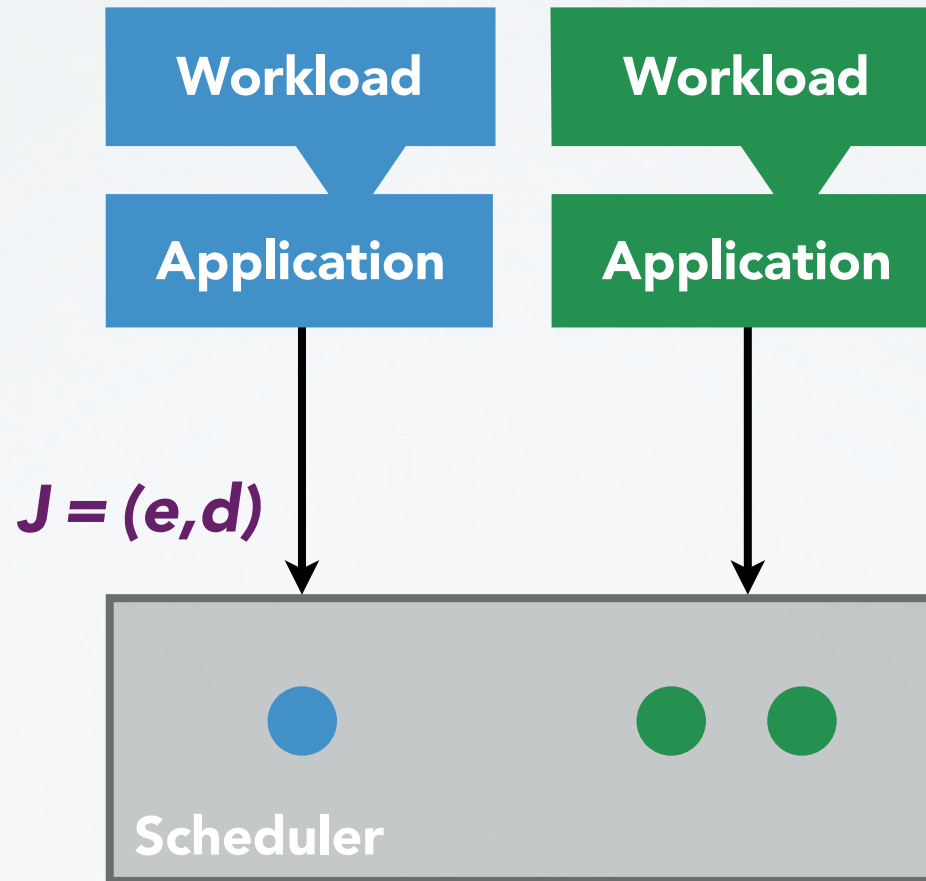
particle density

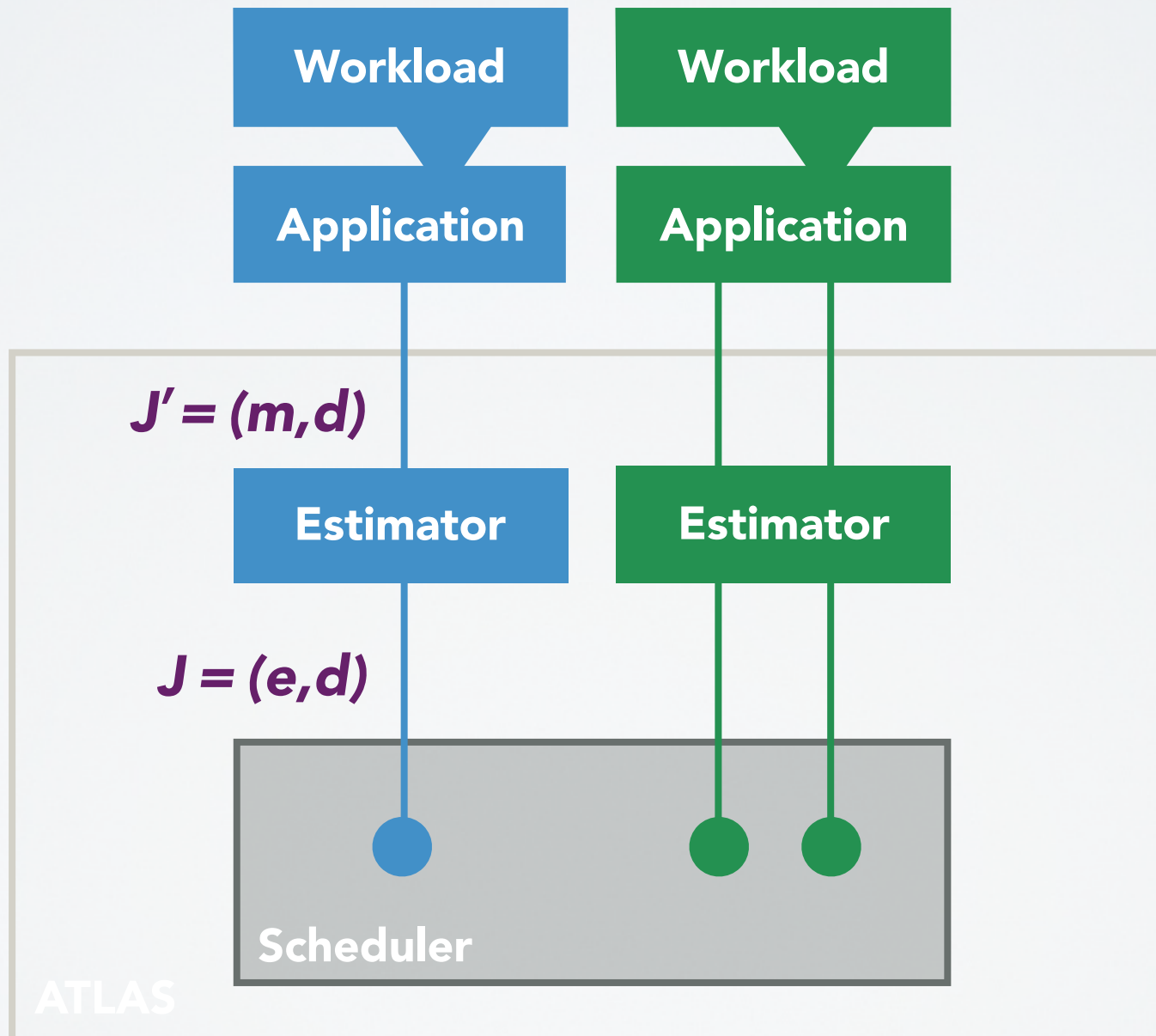
activity centres

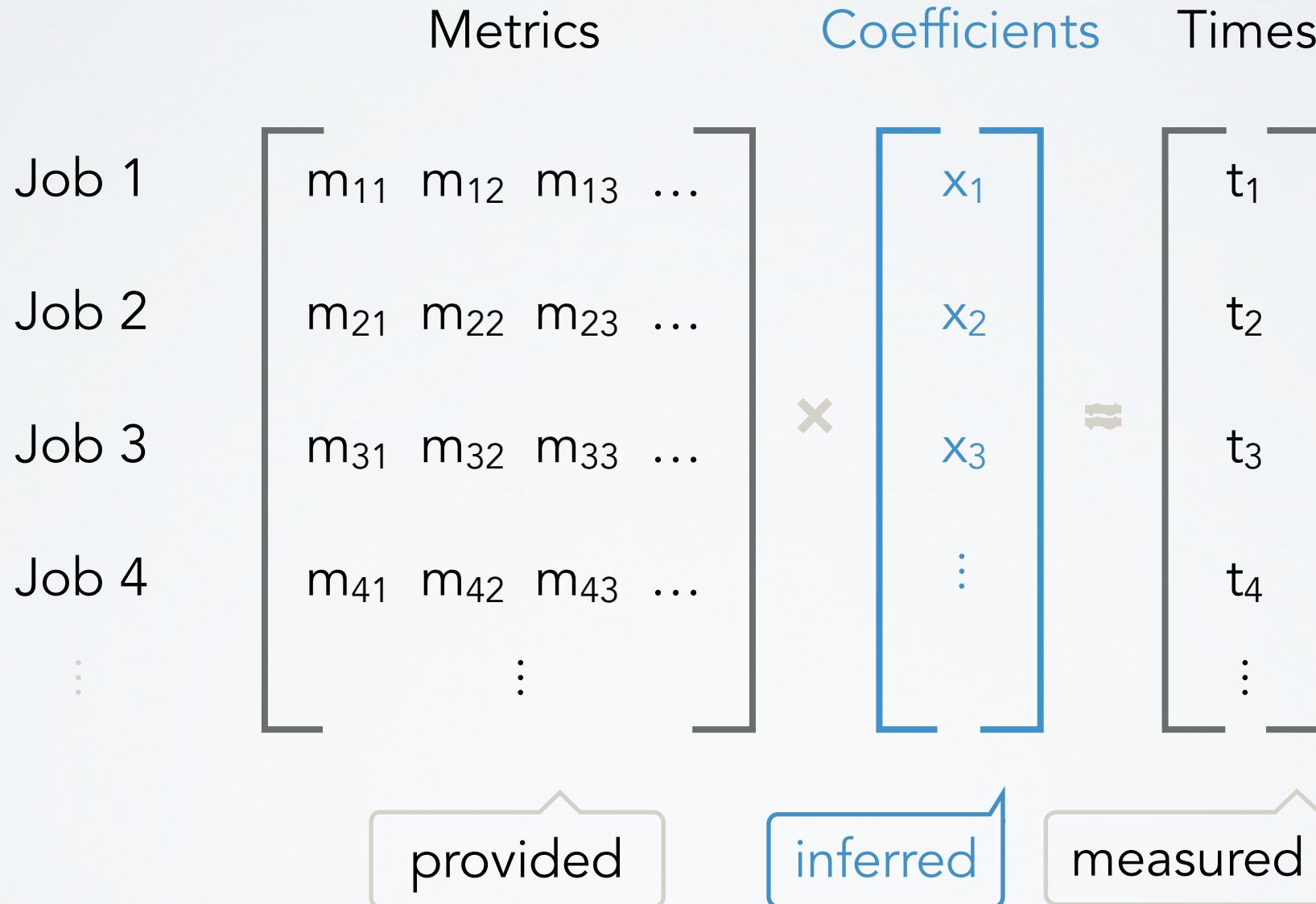
...

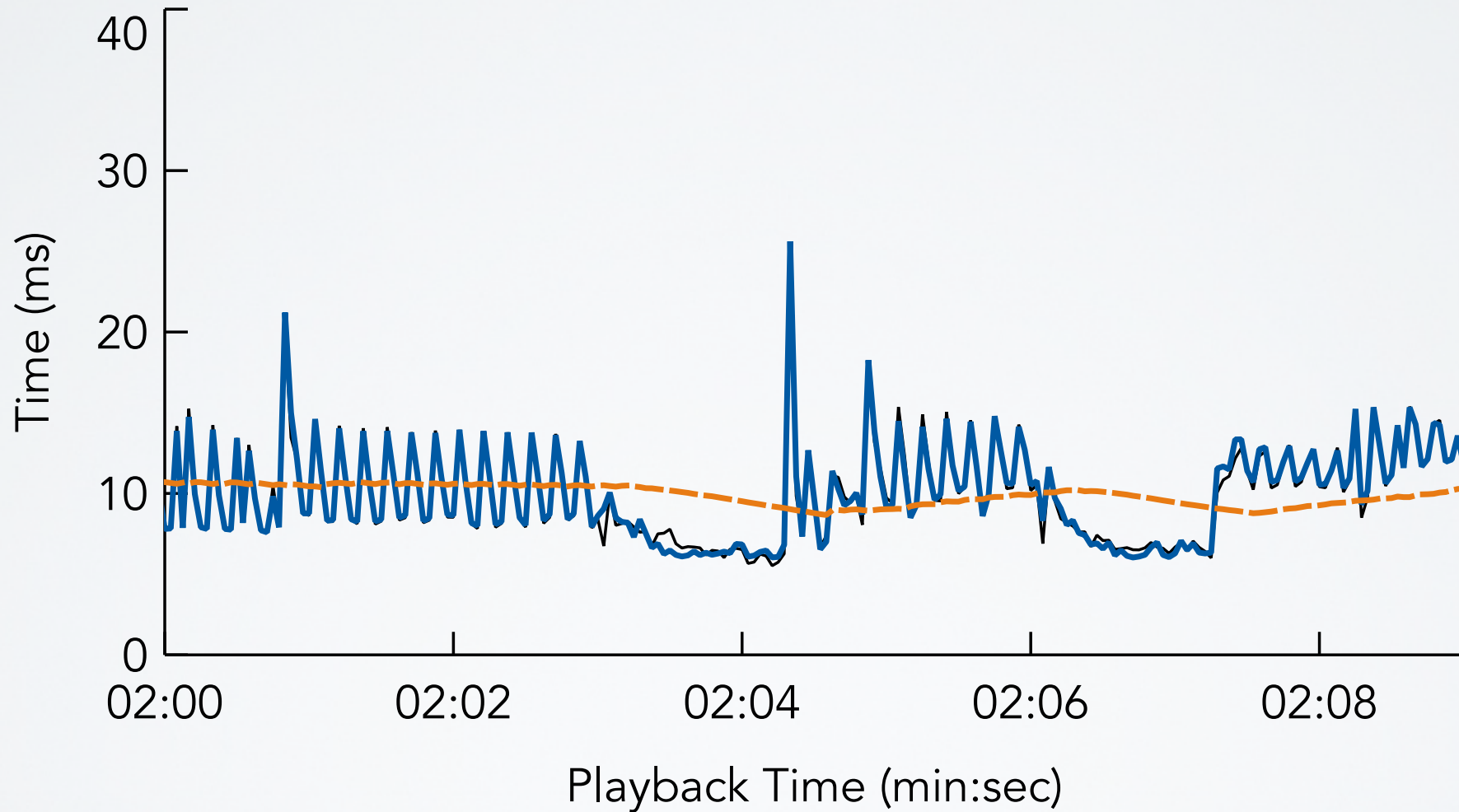


Michael Roitzsch, Practical Real-Time with Look-Ahead Scheduling
PhD Dissertation 2013









cloud density
particle density
activity centres
....

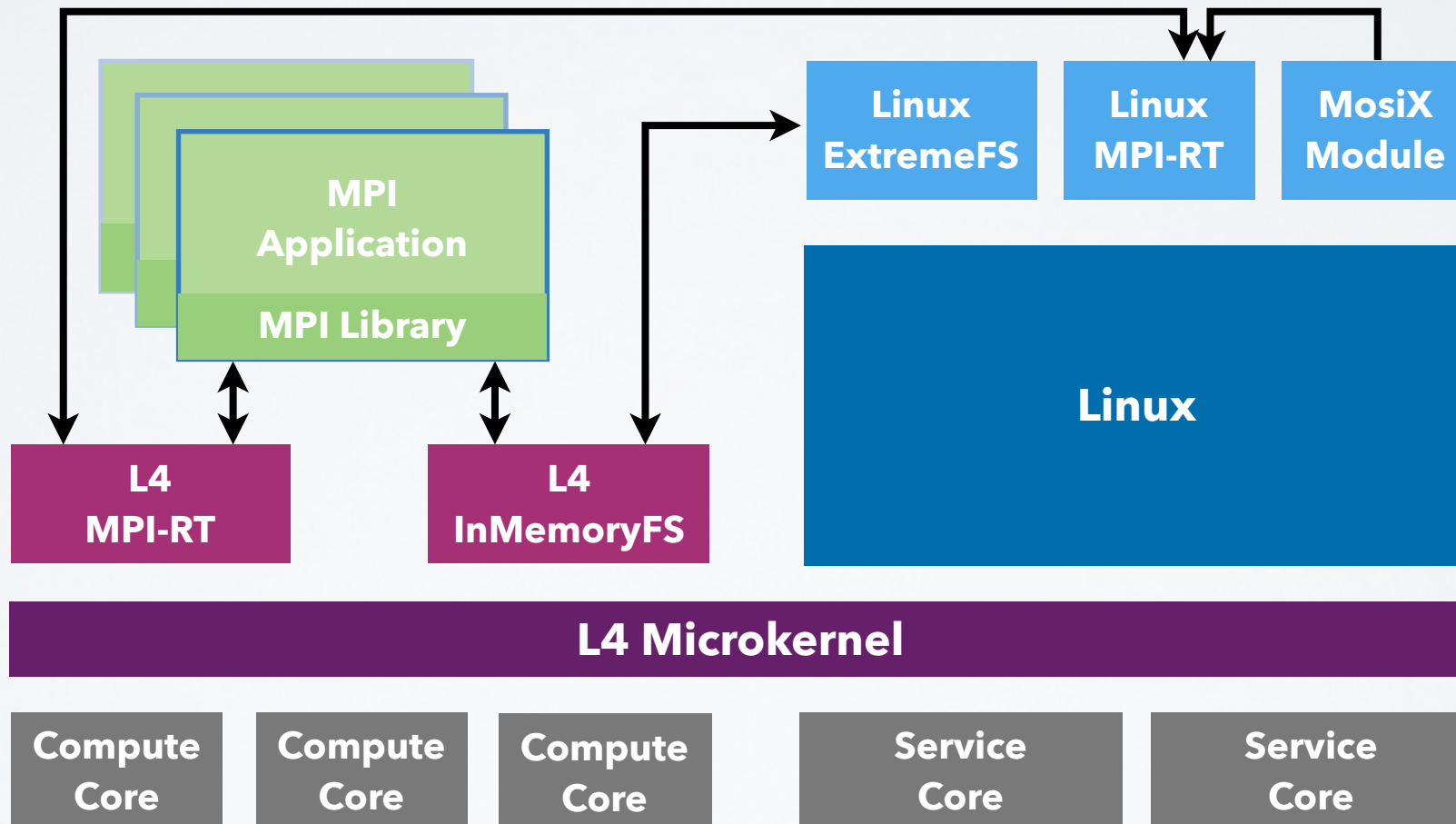


cache misses
cpu cycles
network traffic
....

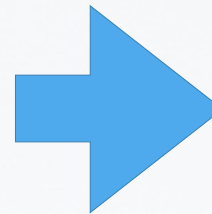
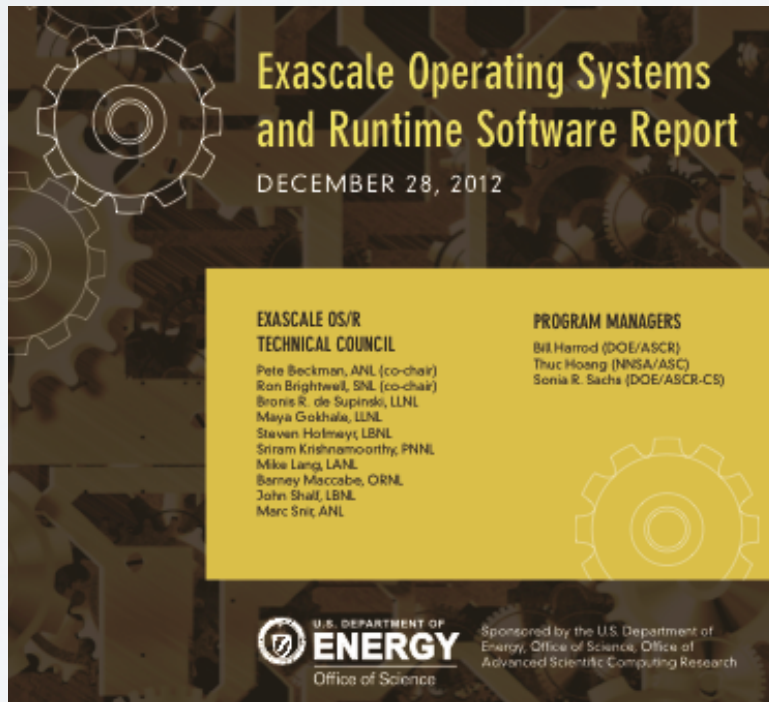


power

Summary



- online decentralized topology- and energy-aware allocation
- programming models interaction with OS
- migration impact on network communication



Fundamental
OS Building Blocks