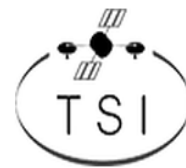


A Scalable Runtime for the ECOSCALE Heterogeneous Exascale Hardware Platform

Paul Harvey

Konstantin Bakanov, Ivor Spence, Dimitrios S. Nikolopoulos







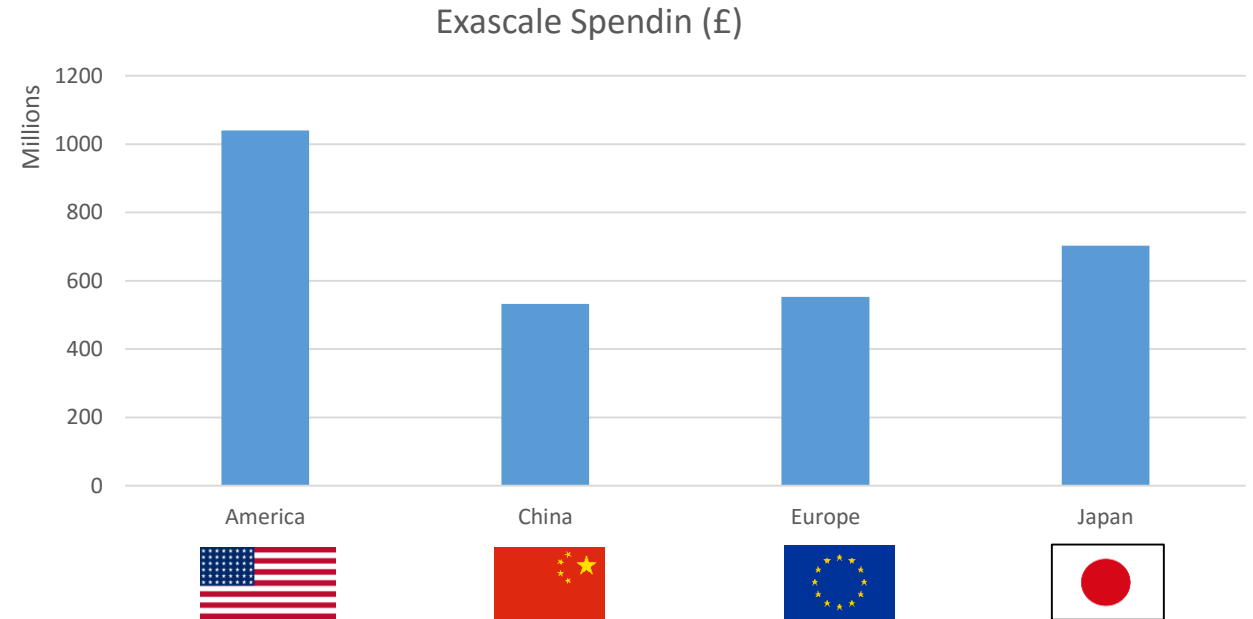
Looking To Discuss and Share Ideas

- No implementation
- No results

- Just design!
 - Intro & Context
 - Hardware
 - Language
 - Runtime Architecture

Exascale: Money

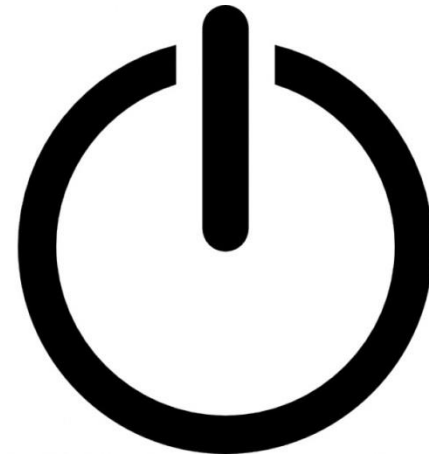
- America : ~\$1500 Million
- Europe : €700 million
- China : 5000 million CNY
- Japan : 110 Billion JPY



Exascale: Brains



Exascale: Problems

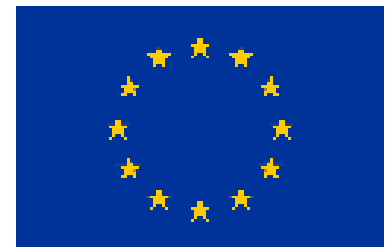


Exascale: Problems



Ecoscale - ecoscale.eu

- Funded till October 2018
- ~£4,000,000
- Building new Hardware
 - Exascale prototype with FPGA focus
- Queen's University working on Software



FPGA



FFT

BitCoin

Matrix Mul

FPGA: Floating point Intensive Calculation

Platform	Time (ns)	W	Energy/Step (nJ)	Obtained By
HD 4400 (GPU)	3.13	15	46.9	Measurement
GTX 960 (GPU)	0.163	120	19.56	Measurement
Quadro K4200 (GPU)	0.204	105	21.42	Measurement
GTX Titan (GPU)	0.0389	375	14.61	Extrapolation
Virtex 7 (FPGA)	0.315	24.4	7.69	Measurement

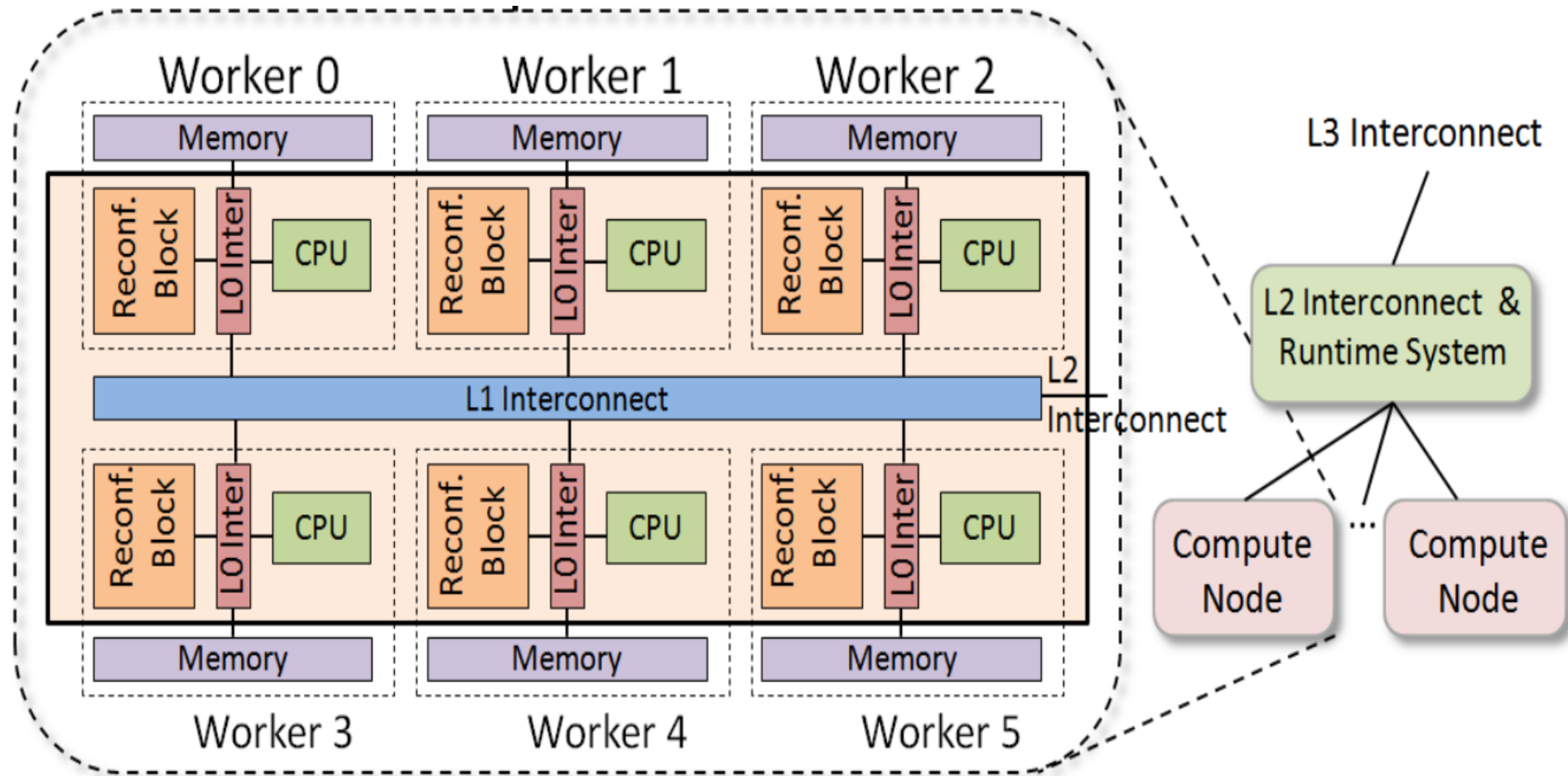
- Compute-intensive, not using global memory
- GPU memory bandwidth is \gg FPGA memory bandwidth
 - GPU DDR4 \sim 8x more than FPGA DDR3

FPGA: Floating point Intensive Calculation

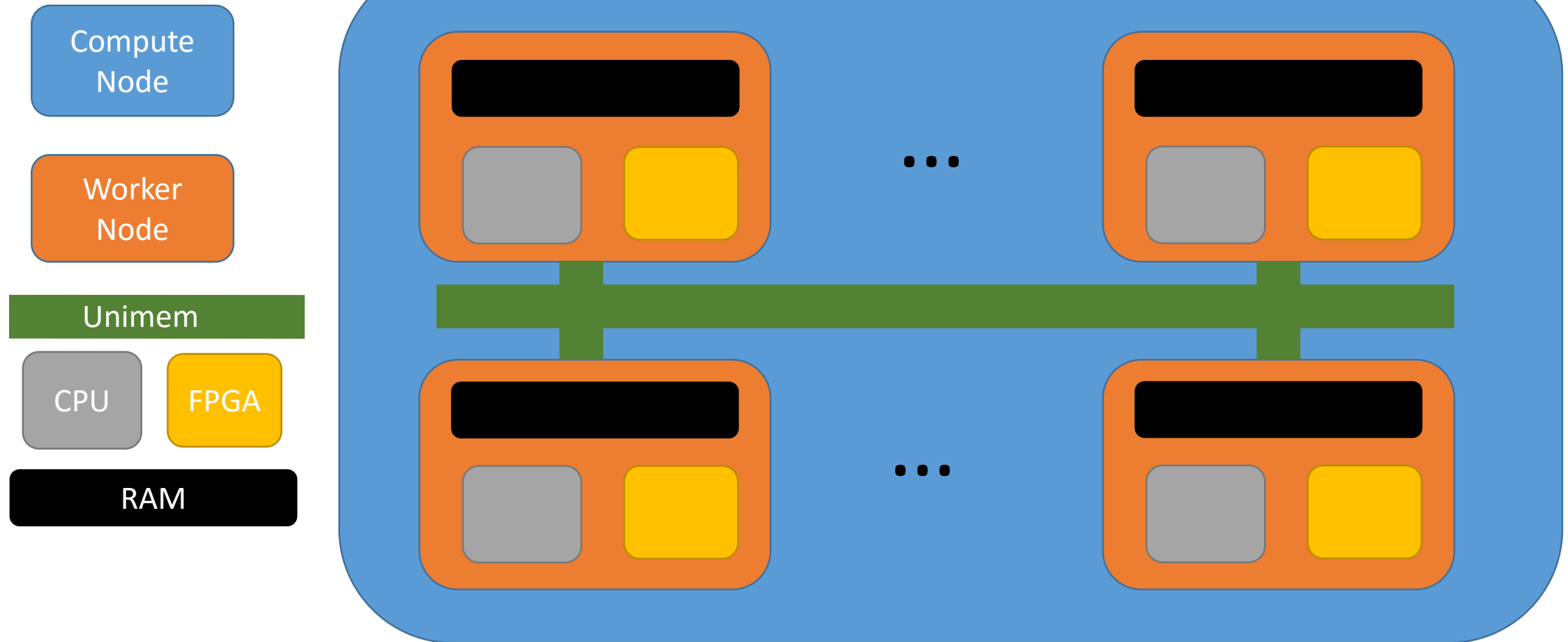
Platform	Time (ns)	W	Energy/Step (nJ)	Obtained By
HD 4400 (GPU)	3.13	15	46.9	Measurement
GTX 960 (GPU)	0.163	120	19.56	Measurement
Quadro K4200 (GPU)	0.204	105	21.42	Measurement
GTX Titan (GPU)	0.0389	375	14.61	Extrapolation
Virtex 7 (FPGA)	0.315	24.4	7.69	Measurement

- Compute-intensive, not using global memory
- GPU memory bandwidth is \gg FPGA memory bandwidth
 - GPU DDR4 \sim 8x more than FPGA DDR3

Architecture



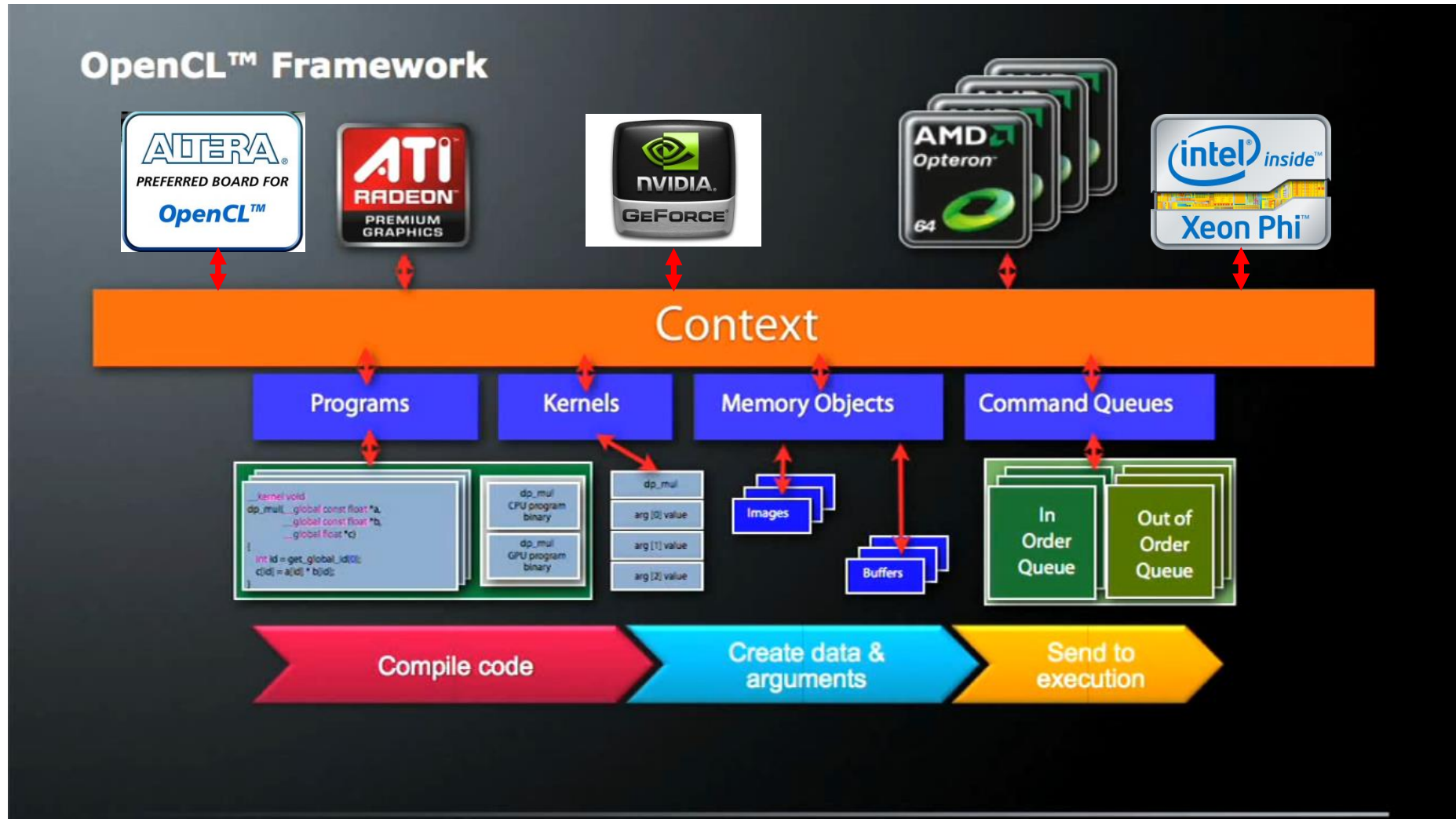
Simplified Architecture



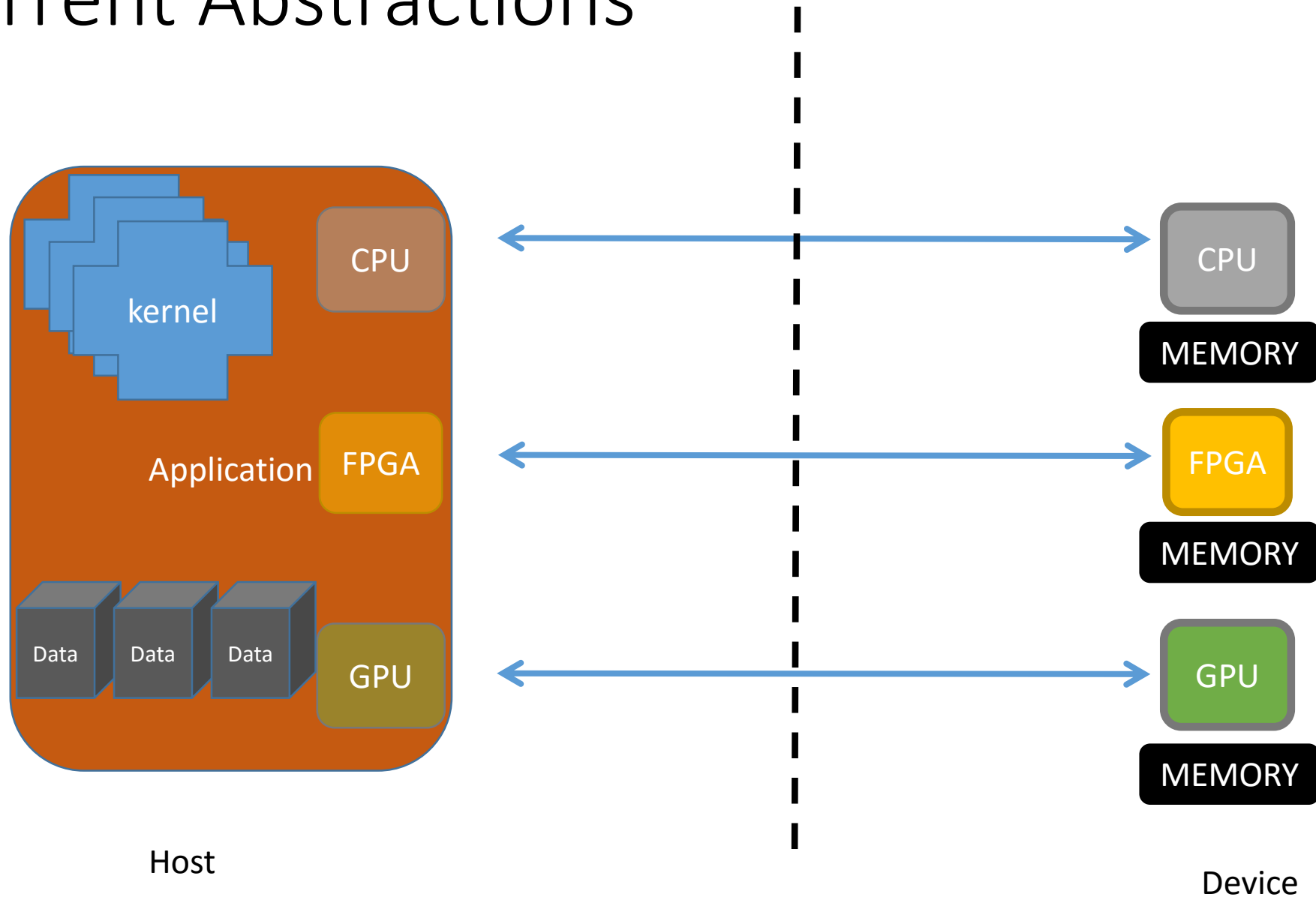
Unimem

- RDMA
- PGAS Address Space
 - One or more single address spaces

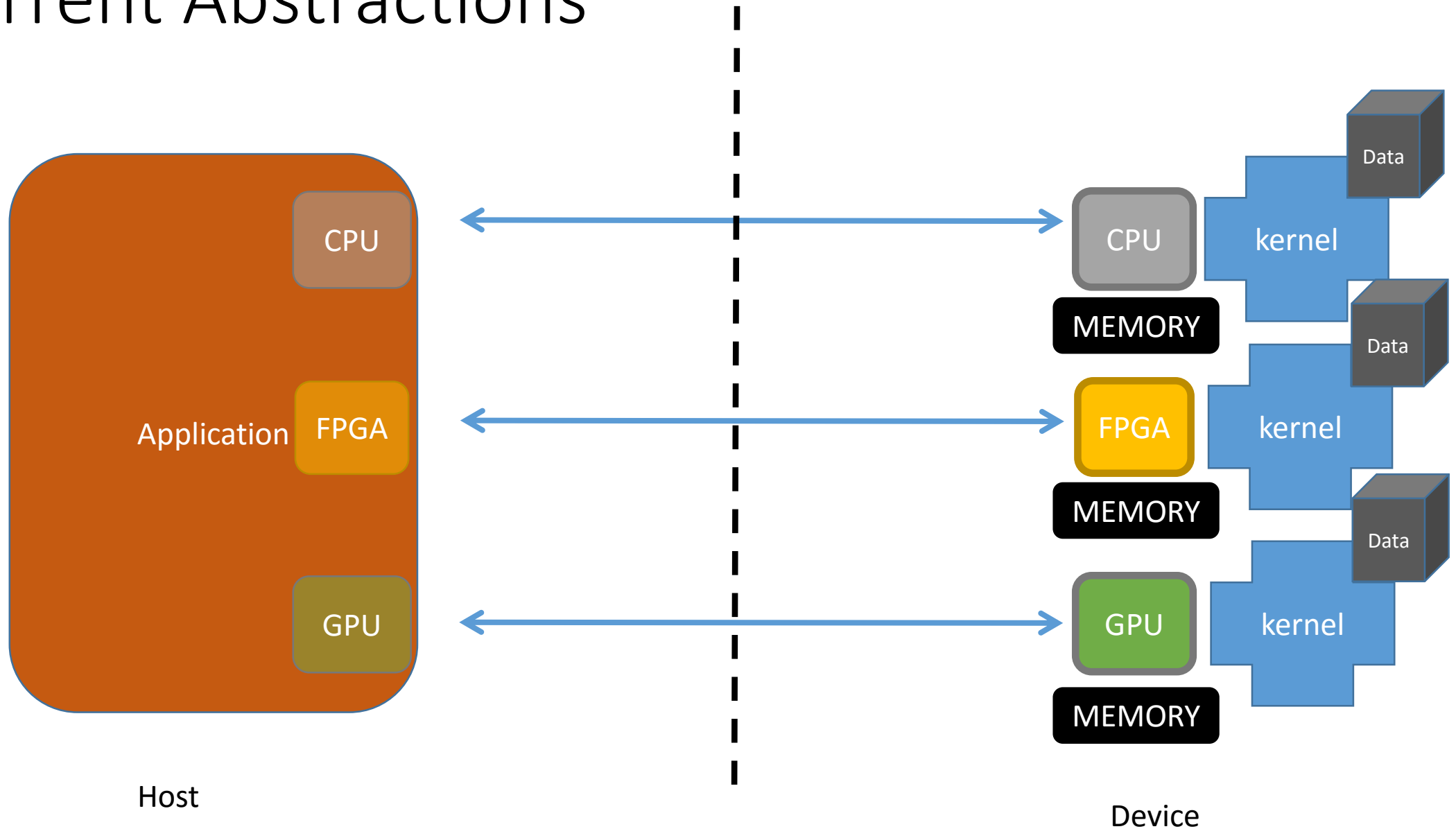
OpenCL



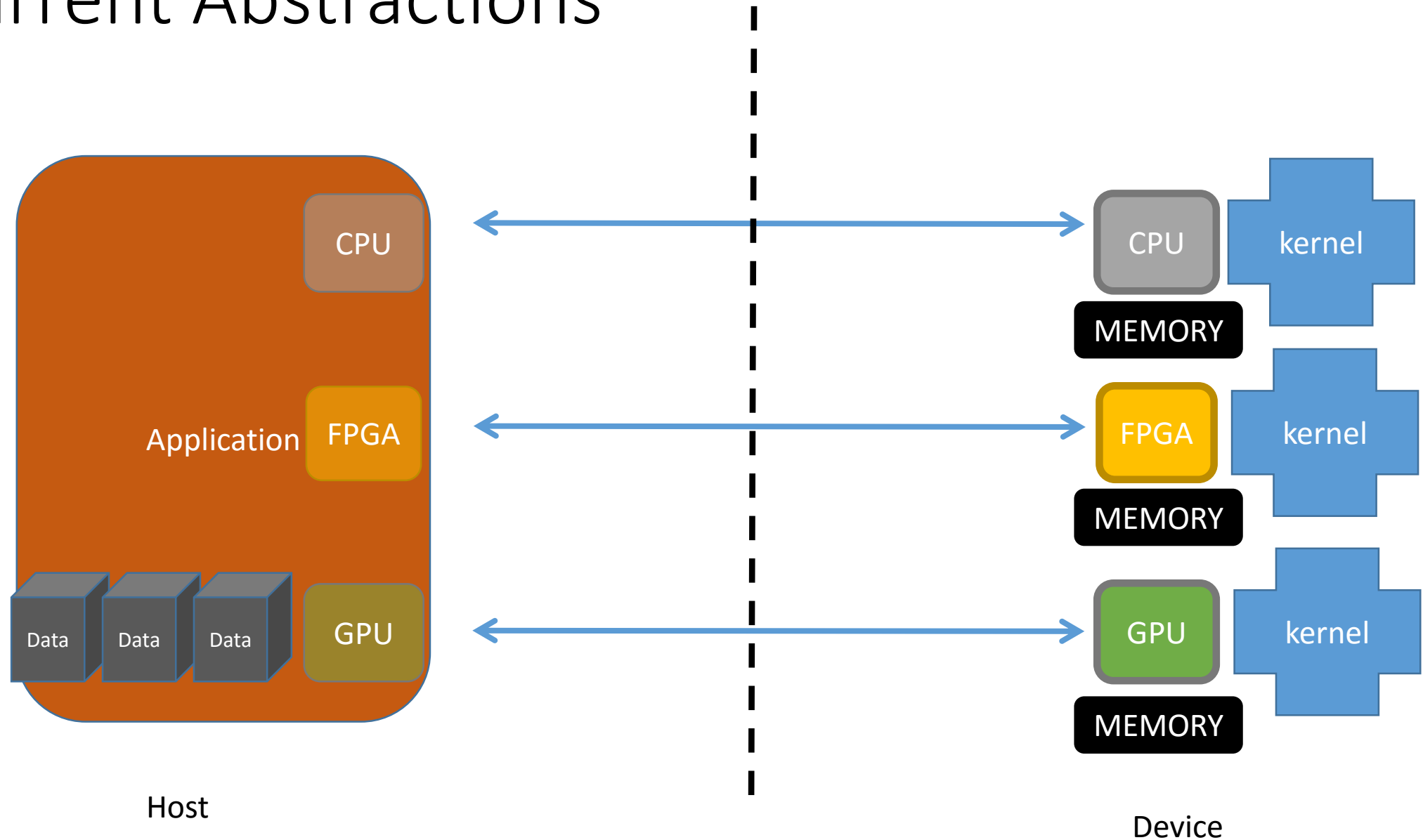
Current Abstractions



Current Abstractions



Current Abstractions



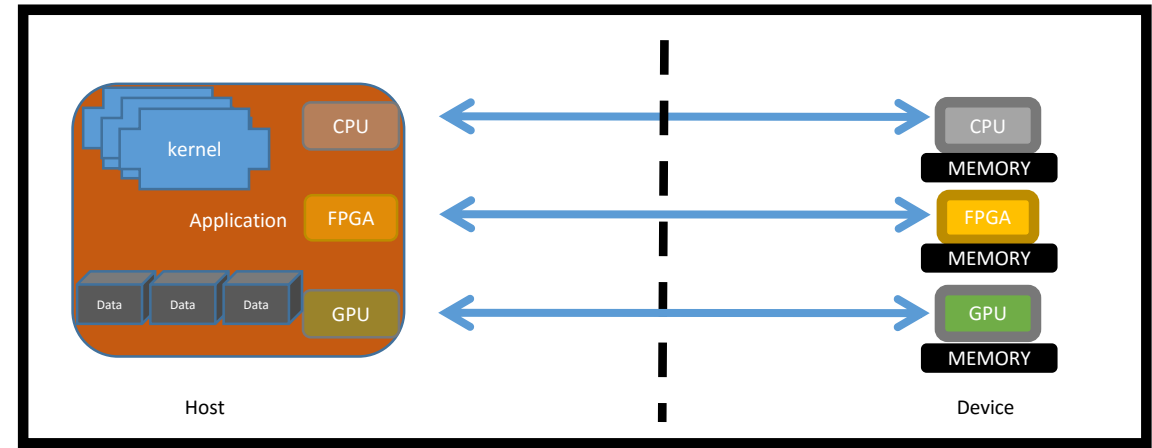
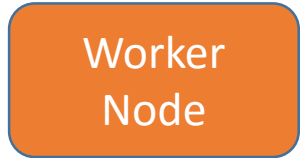
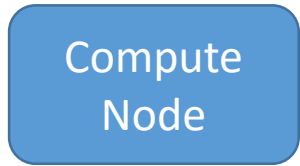
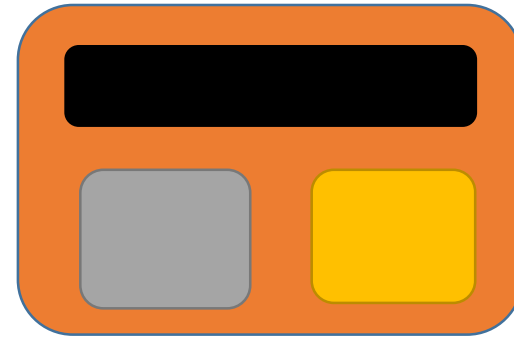
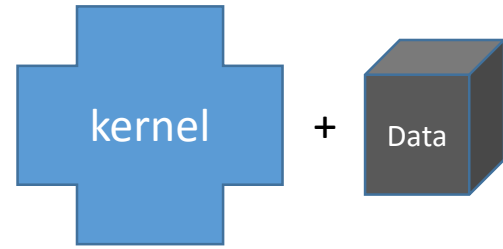
OpenCL

- Simple model
- Widely used in non-hpc
- Standardised
- Lots of activity
 - Industry
 - Academia
- Non-proprietary

Extensions

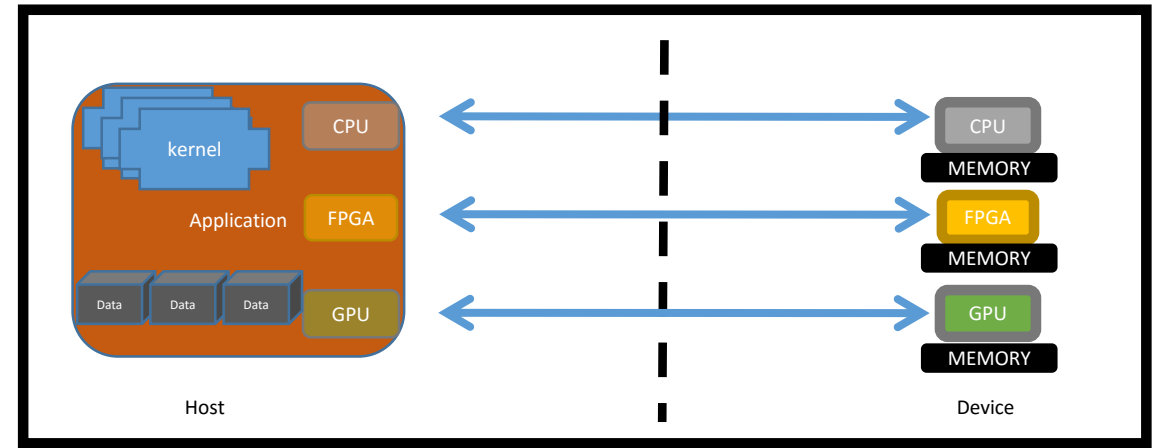
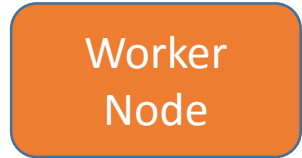
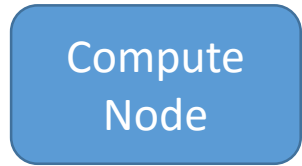
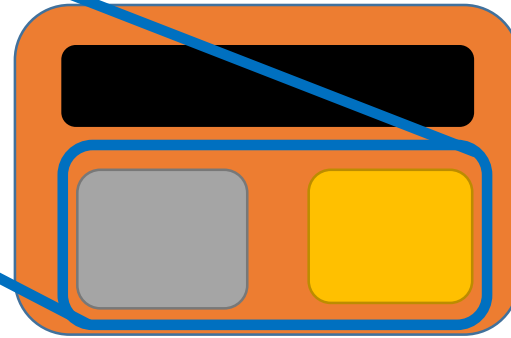
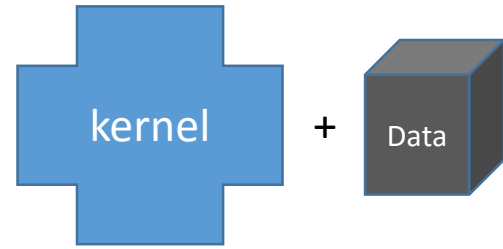
1. New abstractions of multiple hardware devices
 1. Enables scheduler to dynamically go after performance or power
2. New fundamental unit of scheduling
 1. Better scaling across multiple compute devices
 2. Enables kernels to run where a single device has insufficient resources

Worker Abstraction



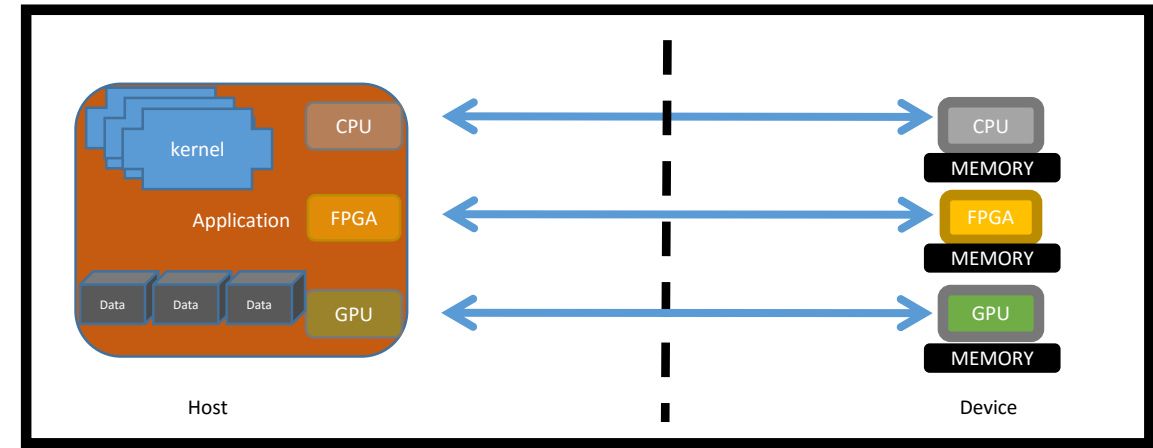
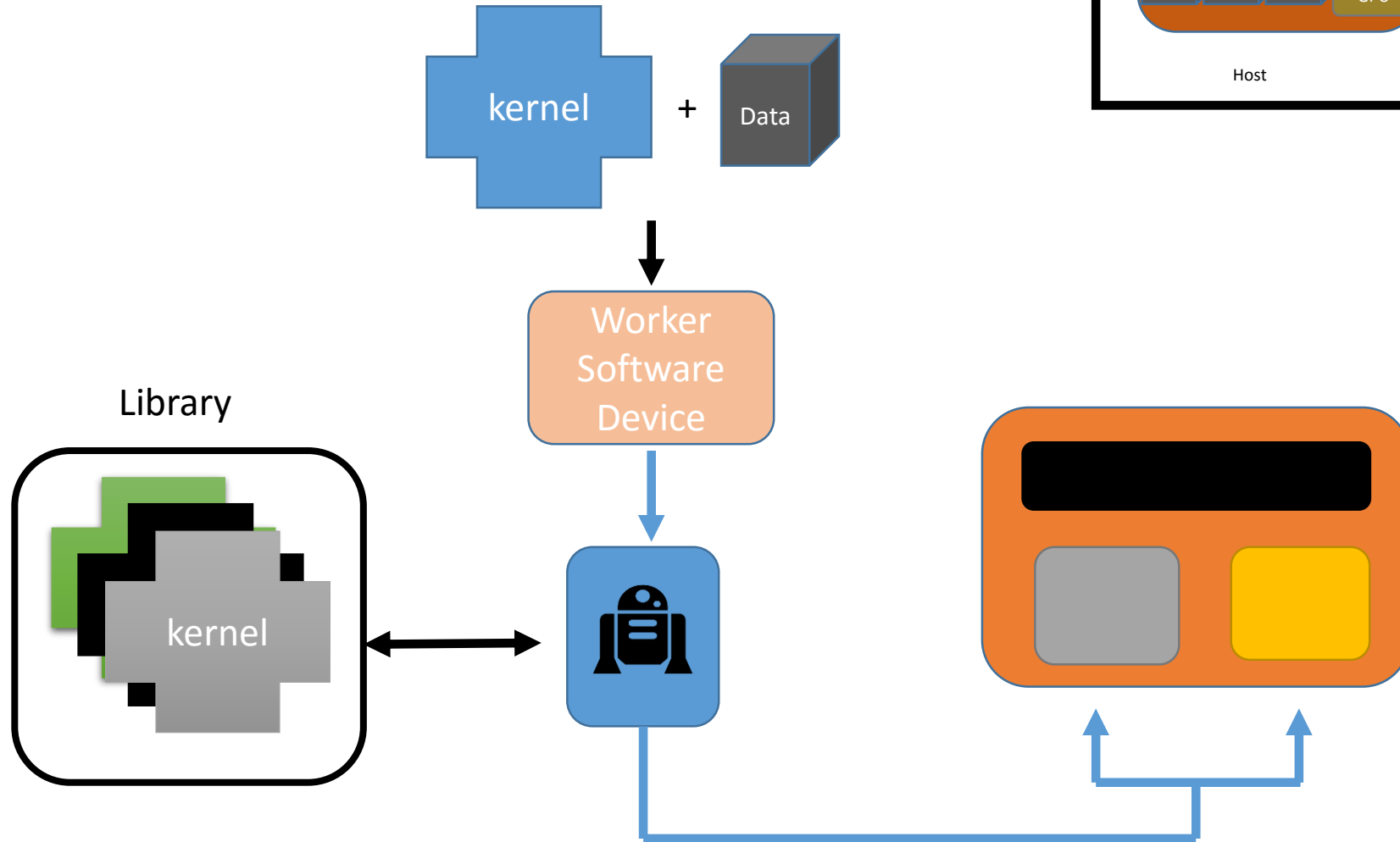
- No change for Programmer
- Scheduler control for power vs. Performance

Worker Abstraction



- No change for Programmer
- Scheduler control for power vs. Performance

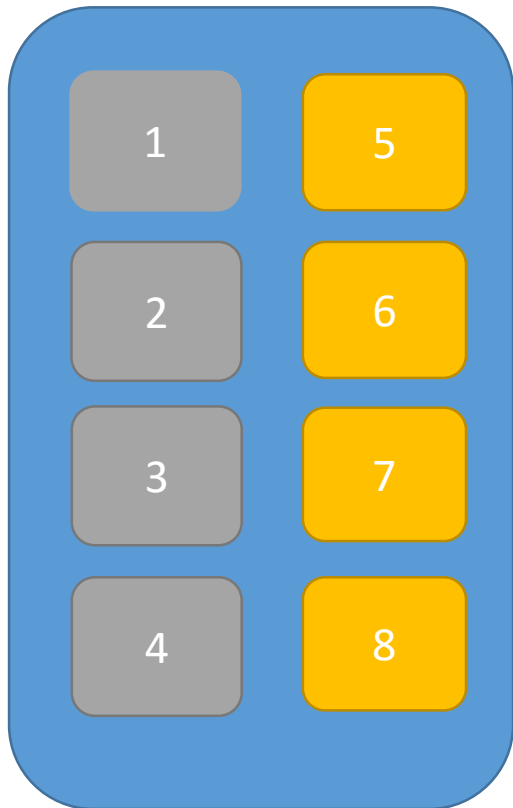
Worker Abstraction



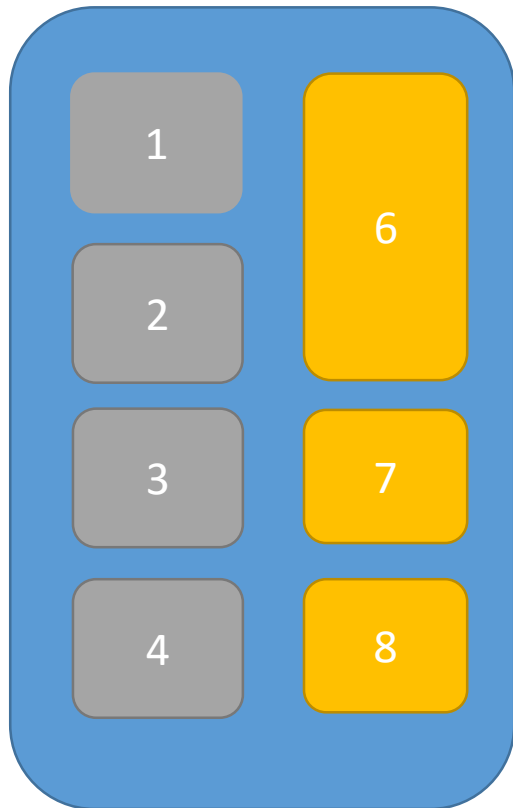
- No change for Programmer
- Scheduler control for power vs. Performance

Abstraction Configurations

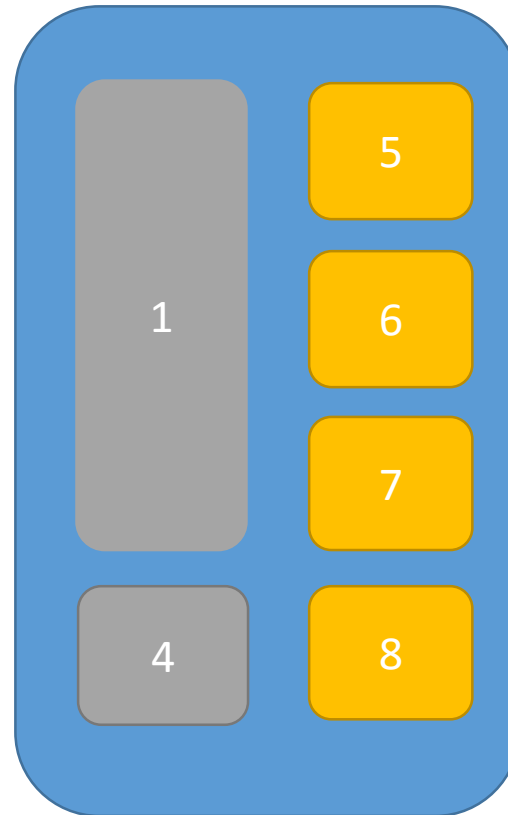
Logical



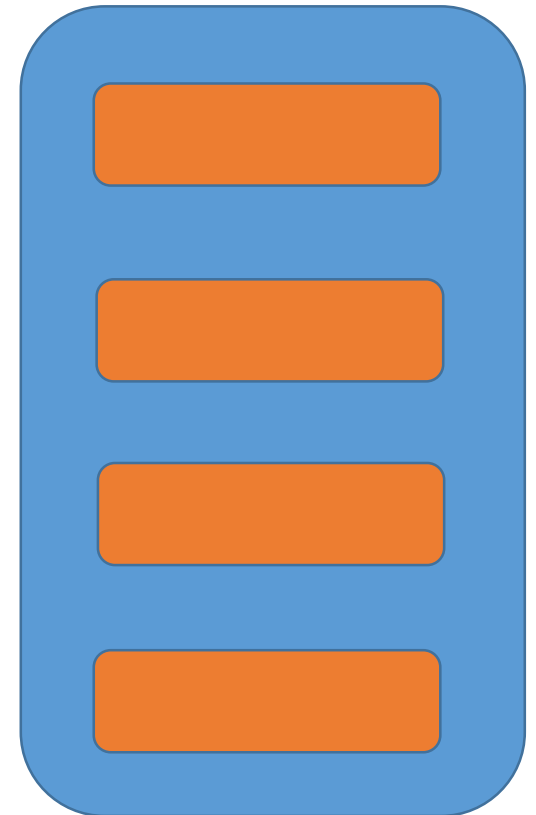
Aggregated FPGA



Aggregated CPU



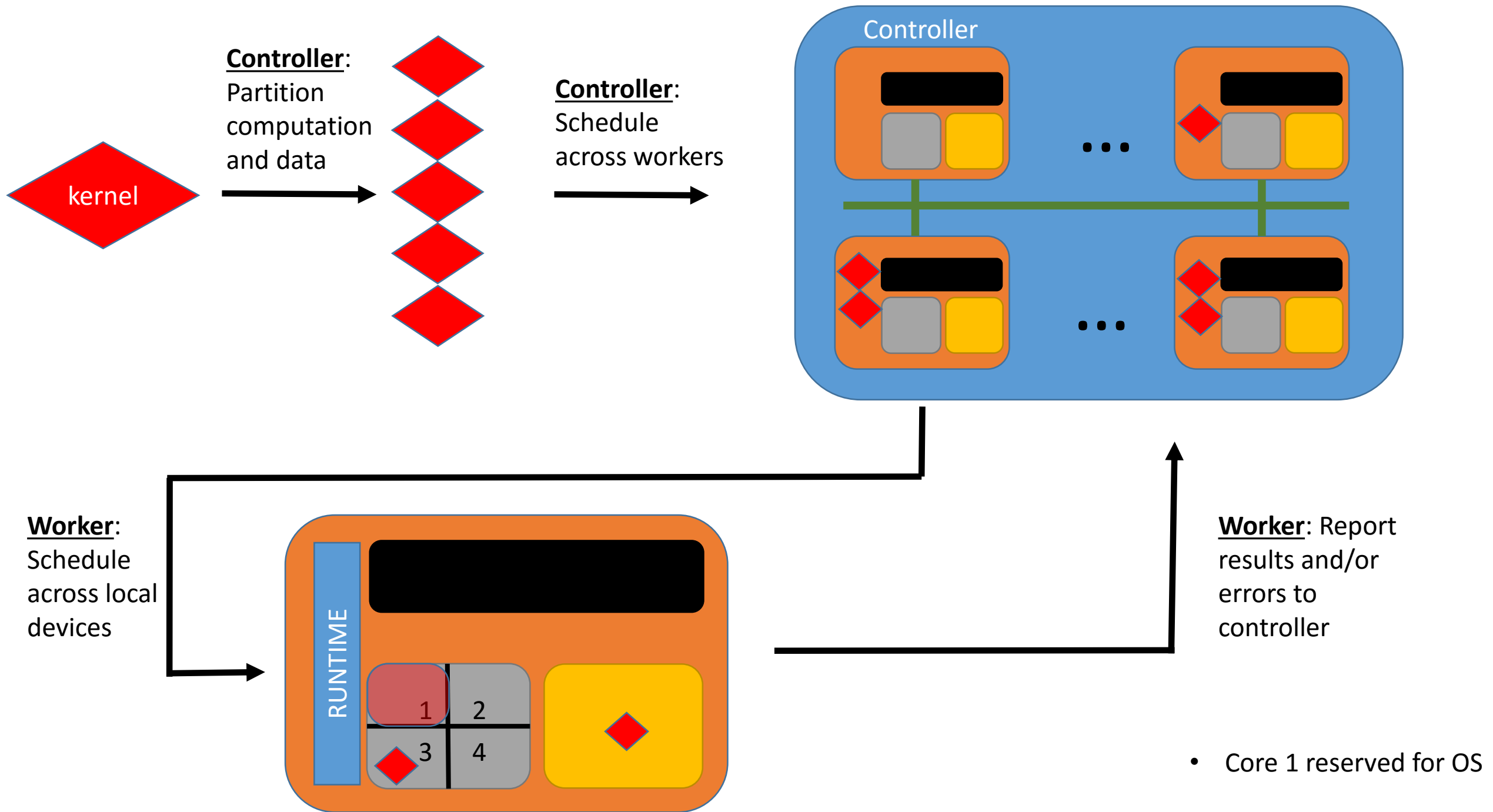
Worker



Scheduling: CPU vs. FPGA



- Machine Learning based on:
 - Runtime performance
 - Kernel input data size
 - CPU/FPGA power consumption
 - Data locality
 - #global memory accesses
 - #branches and loops
- Is a cost model enough?
- How do we determine:
 - a power budget?
 - 100th of current GPU?
 - A performance budget?
 - Current best GPU?



Language – Data Partitioning

```
d_m1 = clCreateBuffer(context,  
    CL_MEM_READ_WRITE,  
    matrix_dim*matrix_dim*sizeof(double),  
    NULL,  
    ecoscale_partition(d_m1, REPLICATE, 0), ←←  
    &errcode);
```

Architecture

Compute Node

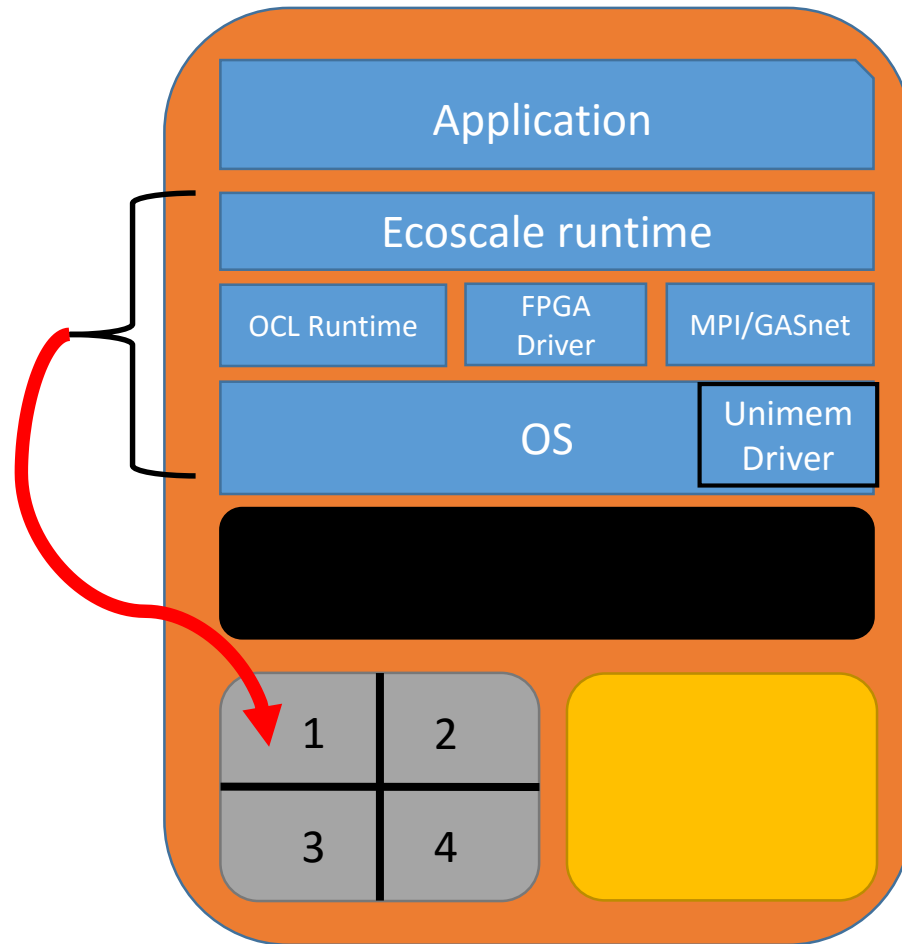
Worker Node

Unimem

CPU

FPGA

RAM



Compute Node

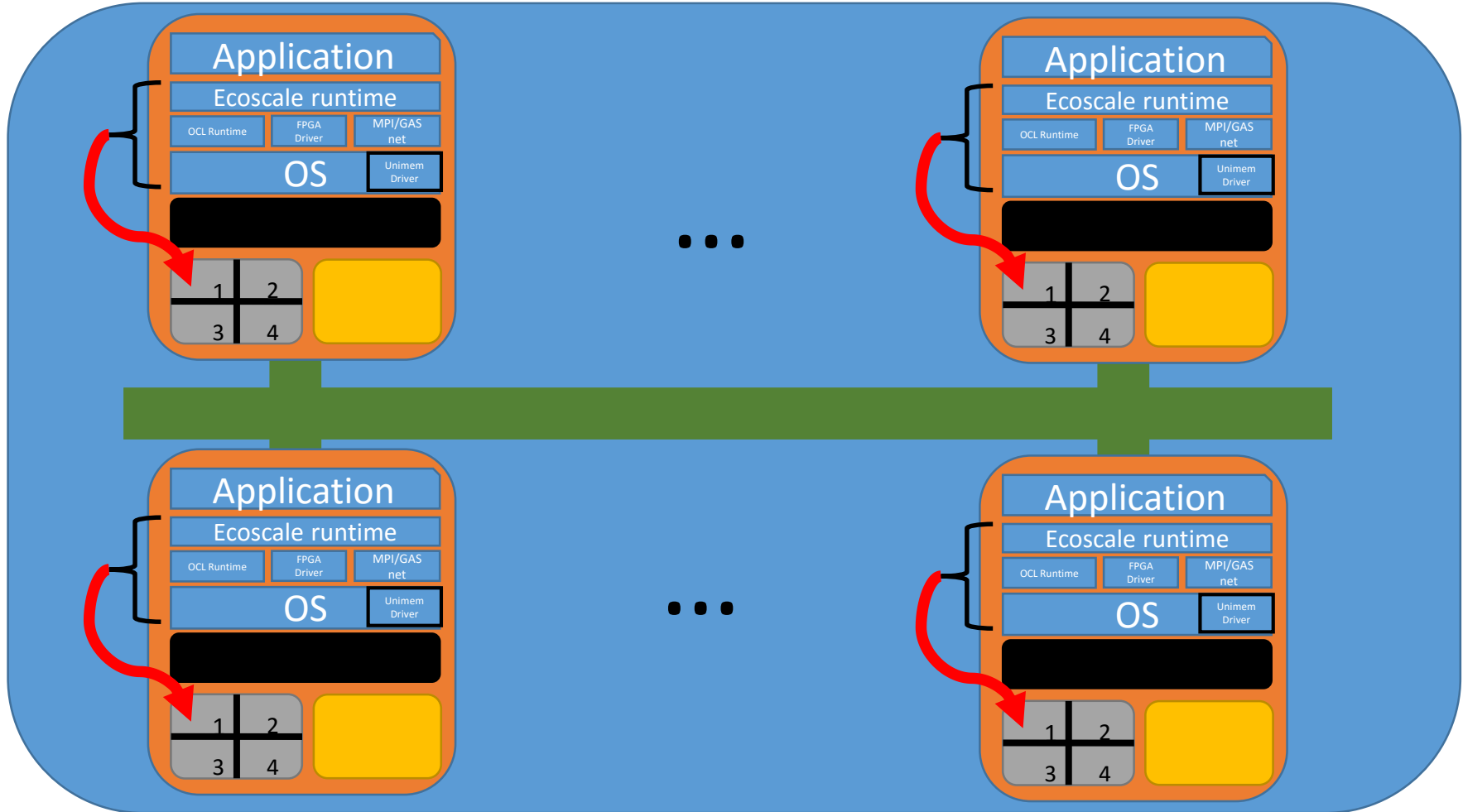
Worker Node

Unimem

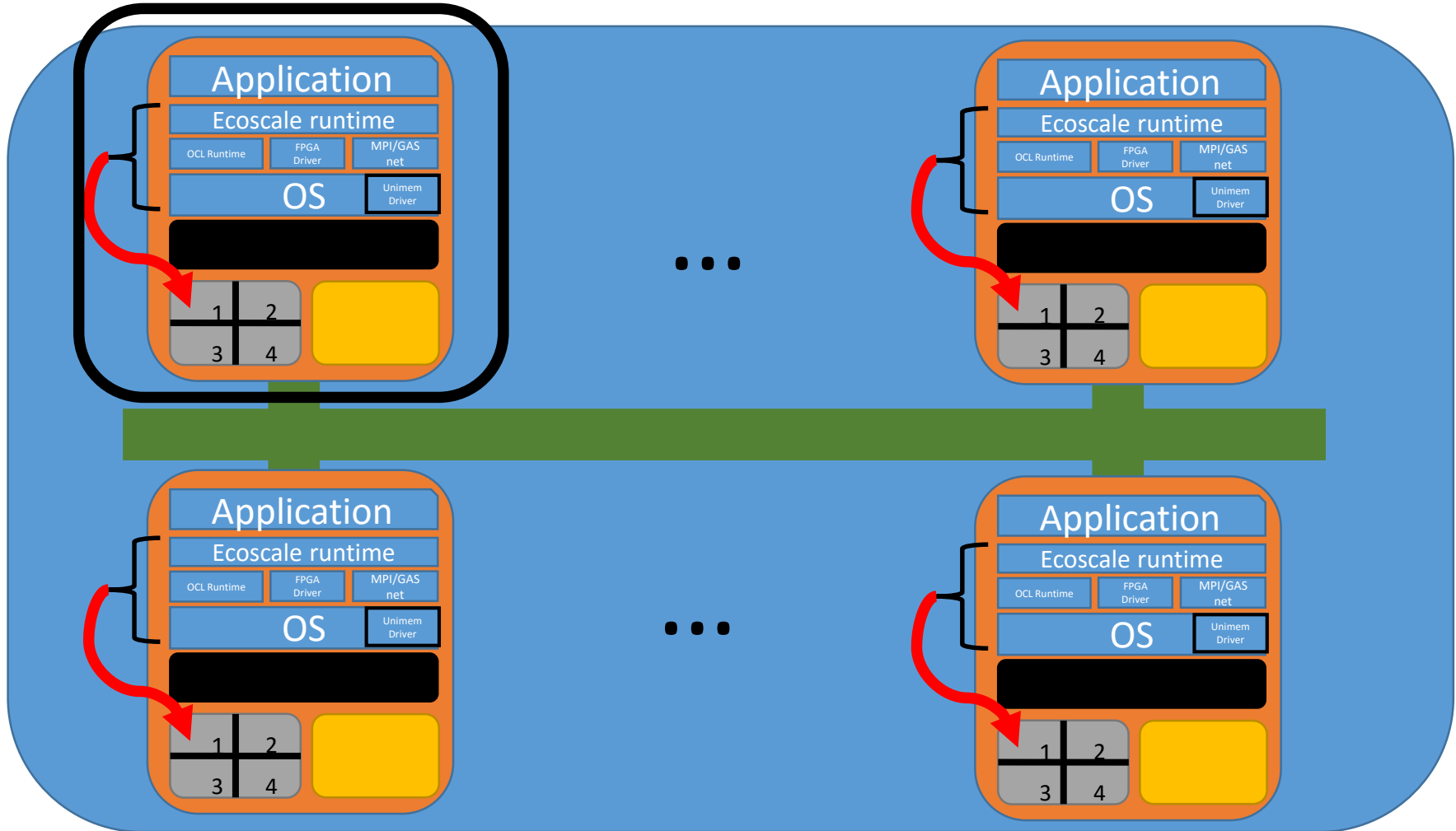
CPU

FPGA

RAM



Controller



Compute Node

Worker Node

Unimem

CPU

FPGA

RAM

Compute Node

Worker Node

Unimem

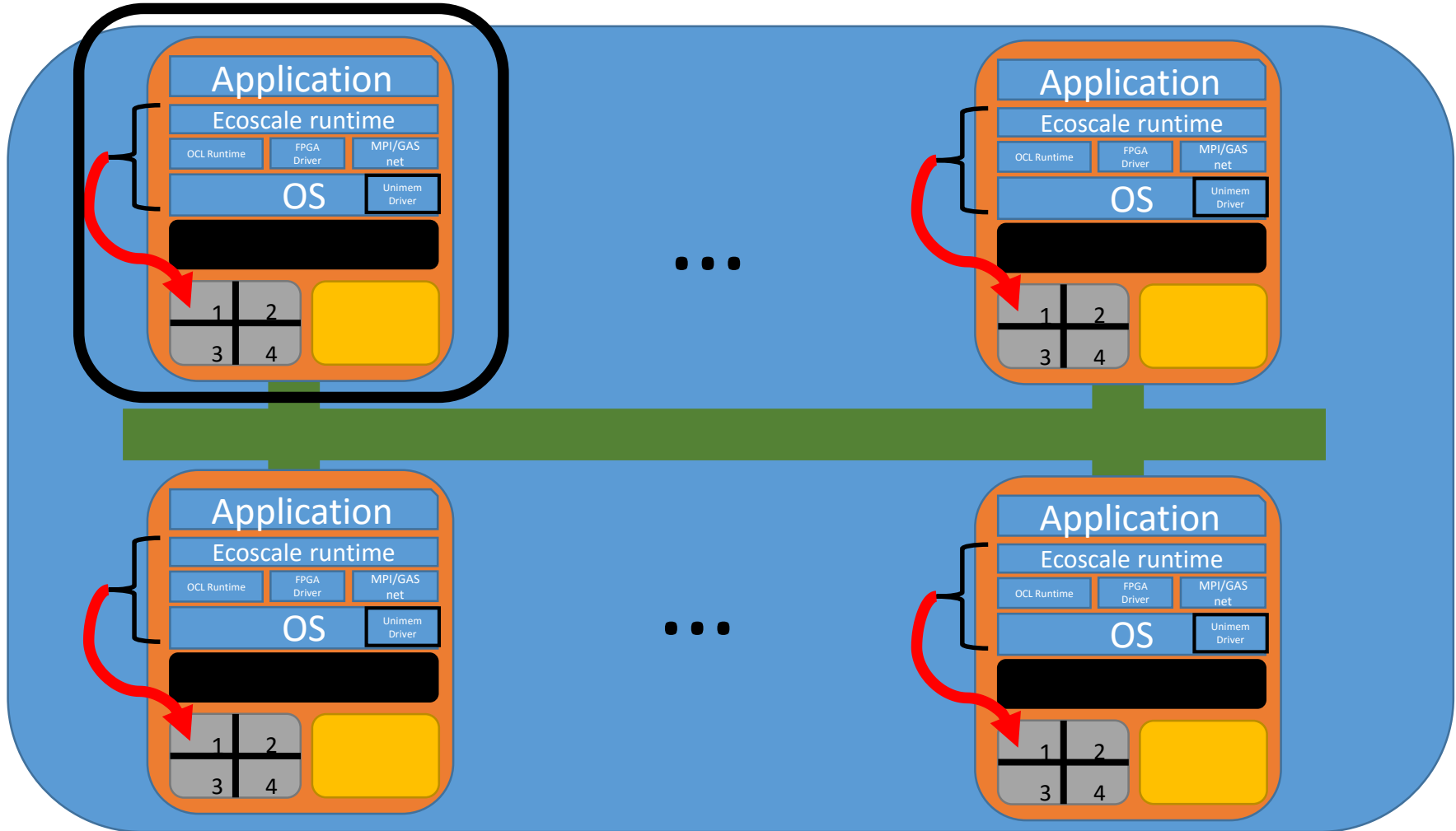
CPU

FPGA

RAM

Controller

Slave



Slave

Slave

Resilience

- Leaders & slaves
- Heartbeats messages
- Checkpointing

Compute Node

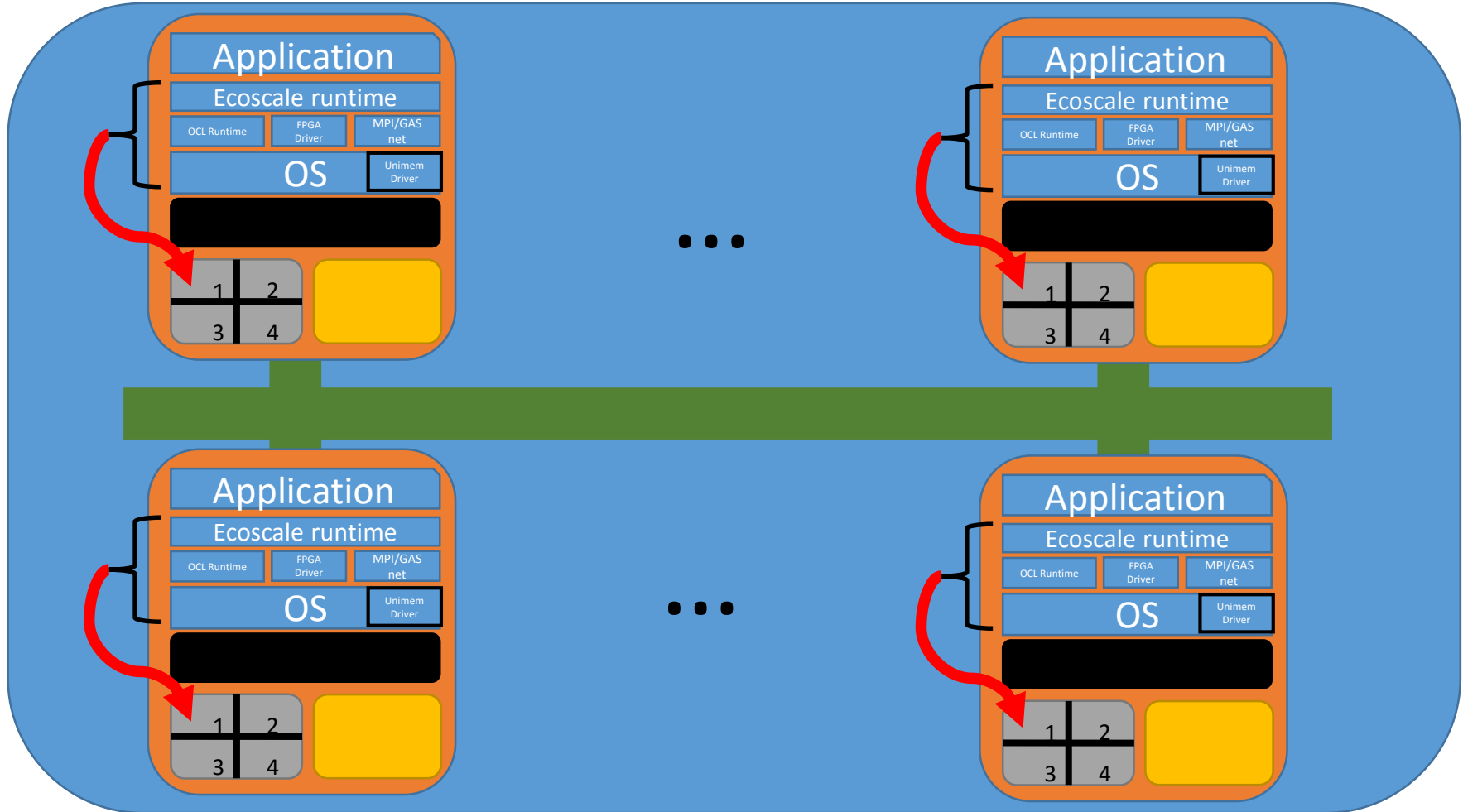
Worker Node

Unimem

CPU

FPGA

RAM



Leadership Election

Compute Node

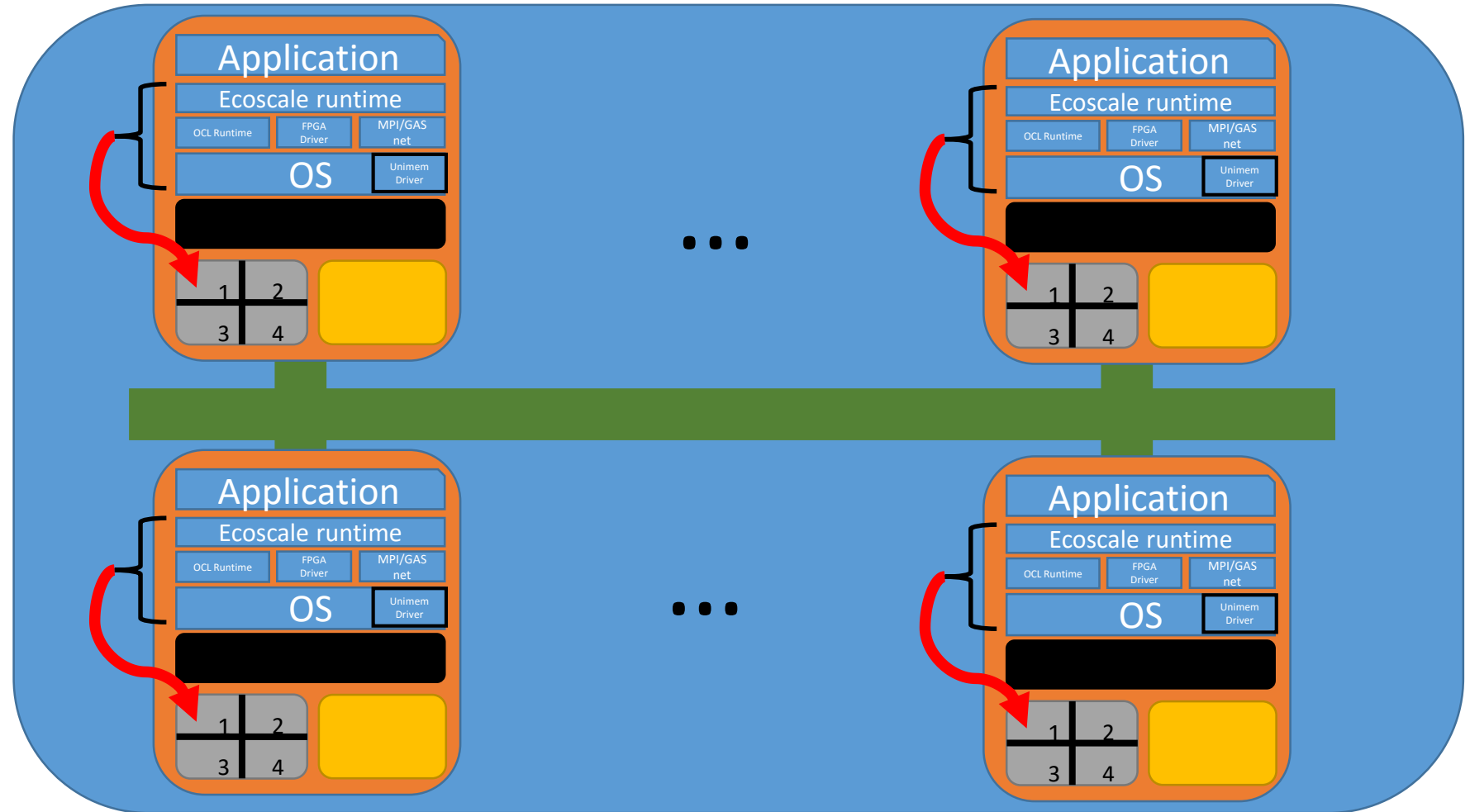
Worker Node

Unimem

CPU

FPGA

RAM



Compute Node

Worker Node

Unimem

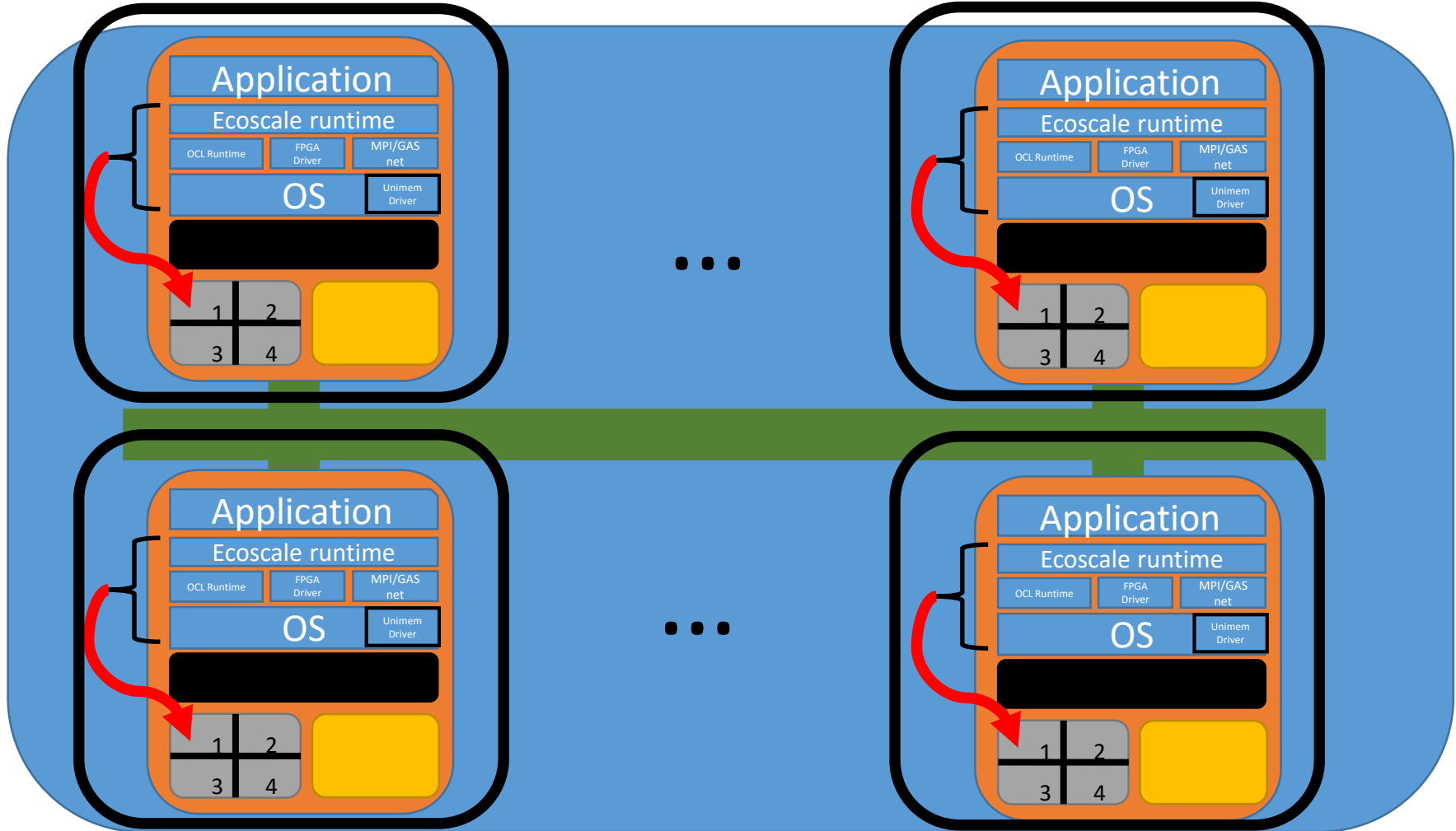
CPU

FPGA

RAM

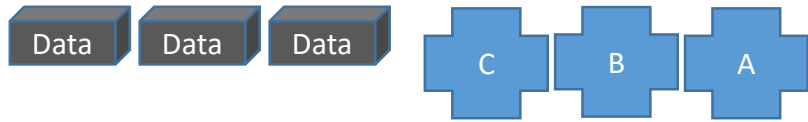
Controller

Slave (Backup)



Slave

Slave



Accounting Log

Controller

Slave (Backup)

Compute Node

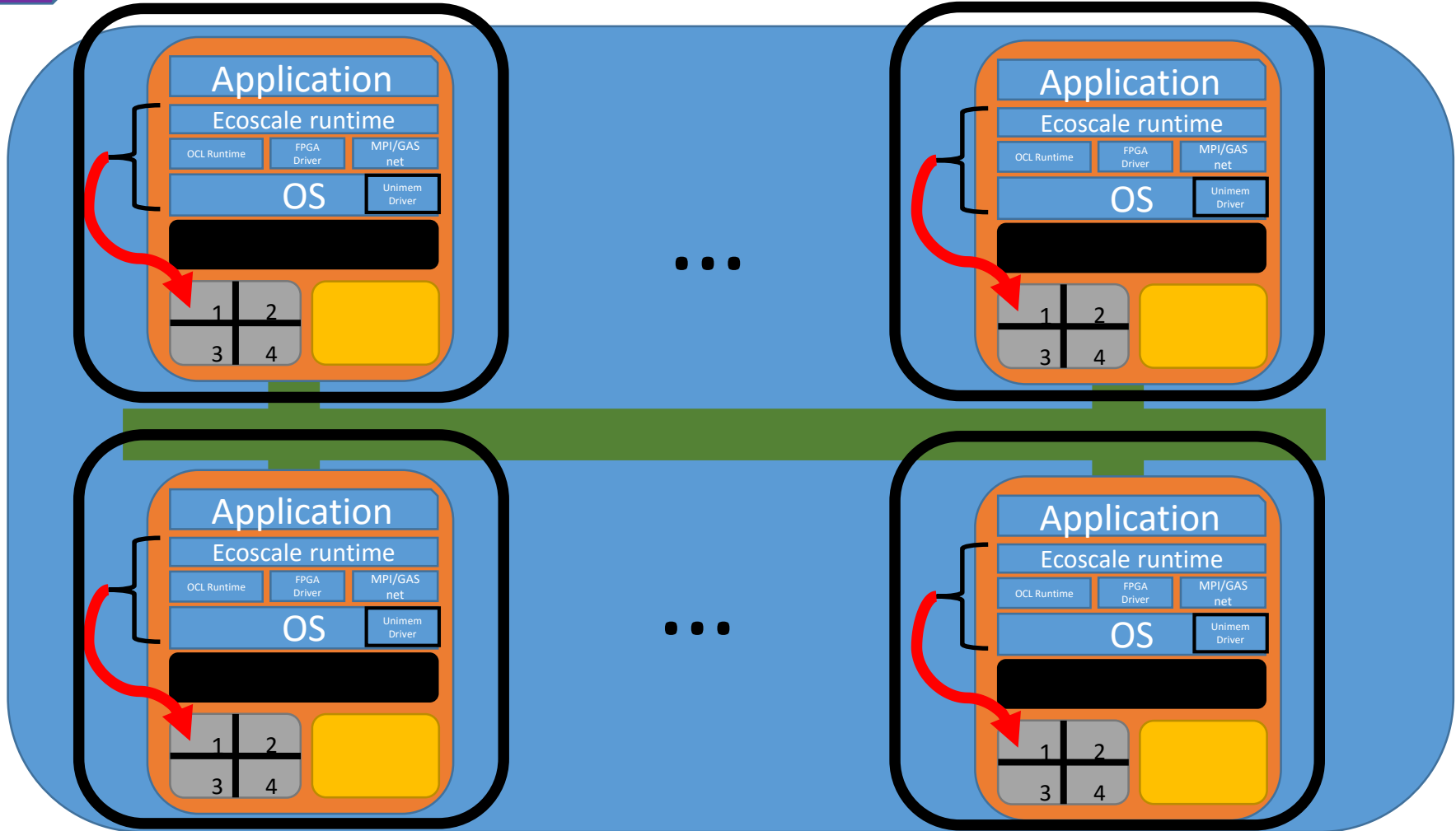
Worker Node

Unimem

CPU

FPGA

RAM



Slave

Slave

Accounting Log

Compute Node

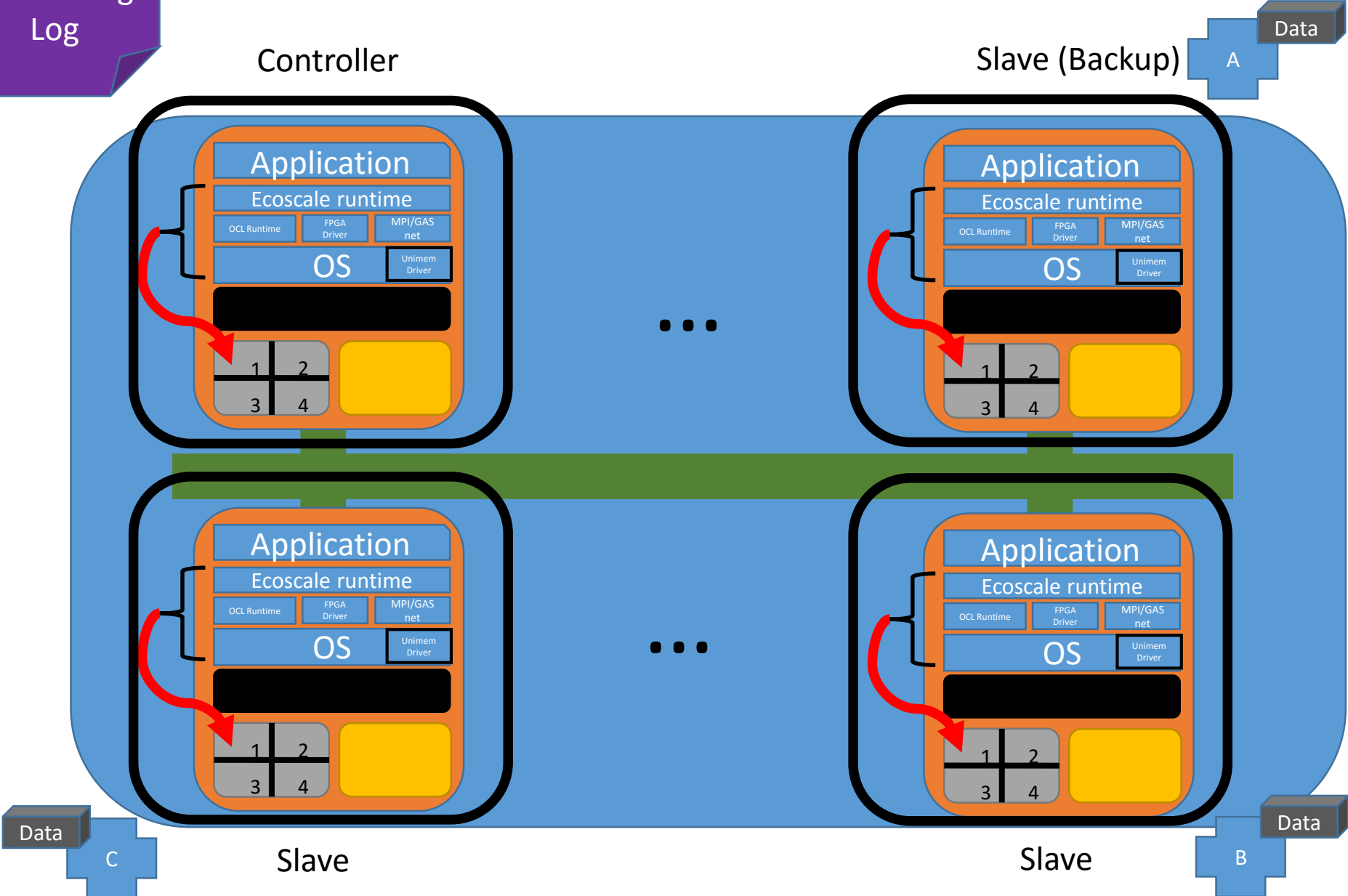
Worker Node

Unimem

CPU

FPGA

RAM



Controller

Slave (Backup)

Slave

Slave

Data

C

Data

A

Data

B

Accounting Log

Compute Node

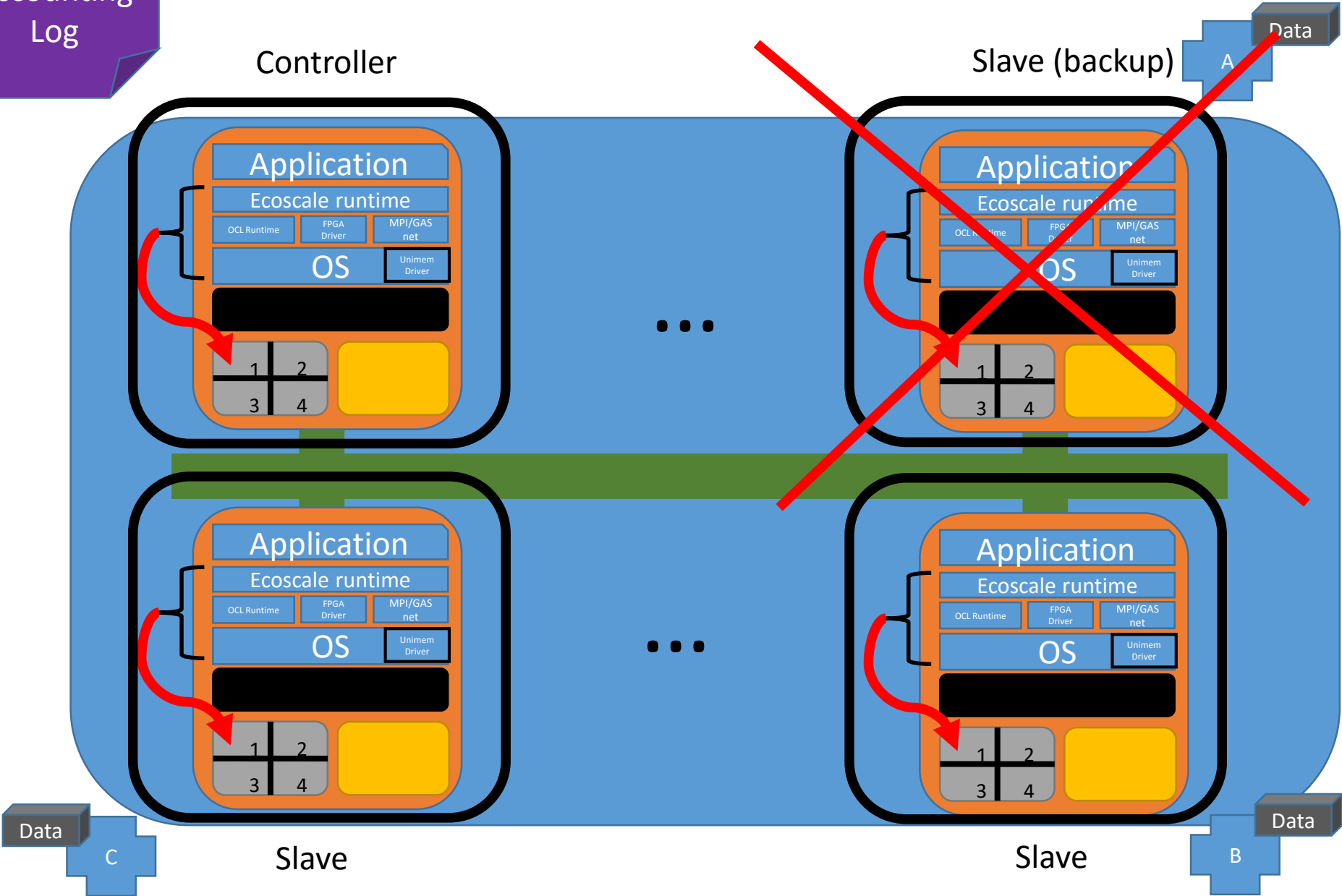
Worker Node

Unimem

CPU

FPGA

RAM



Accounting Log

Controller

Slave (backup)

Data

A

Data

C

Slave

Slave

Data

B

Compute Node

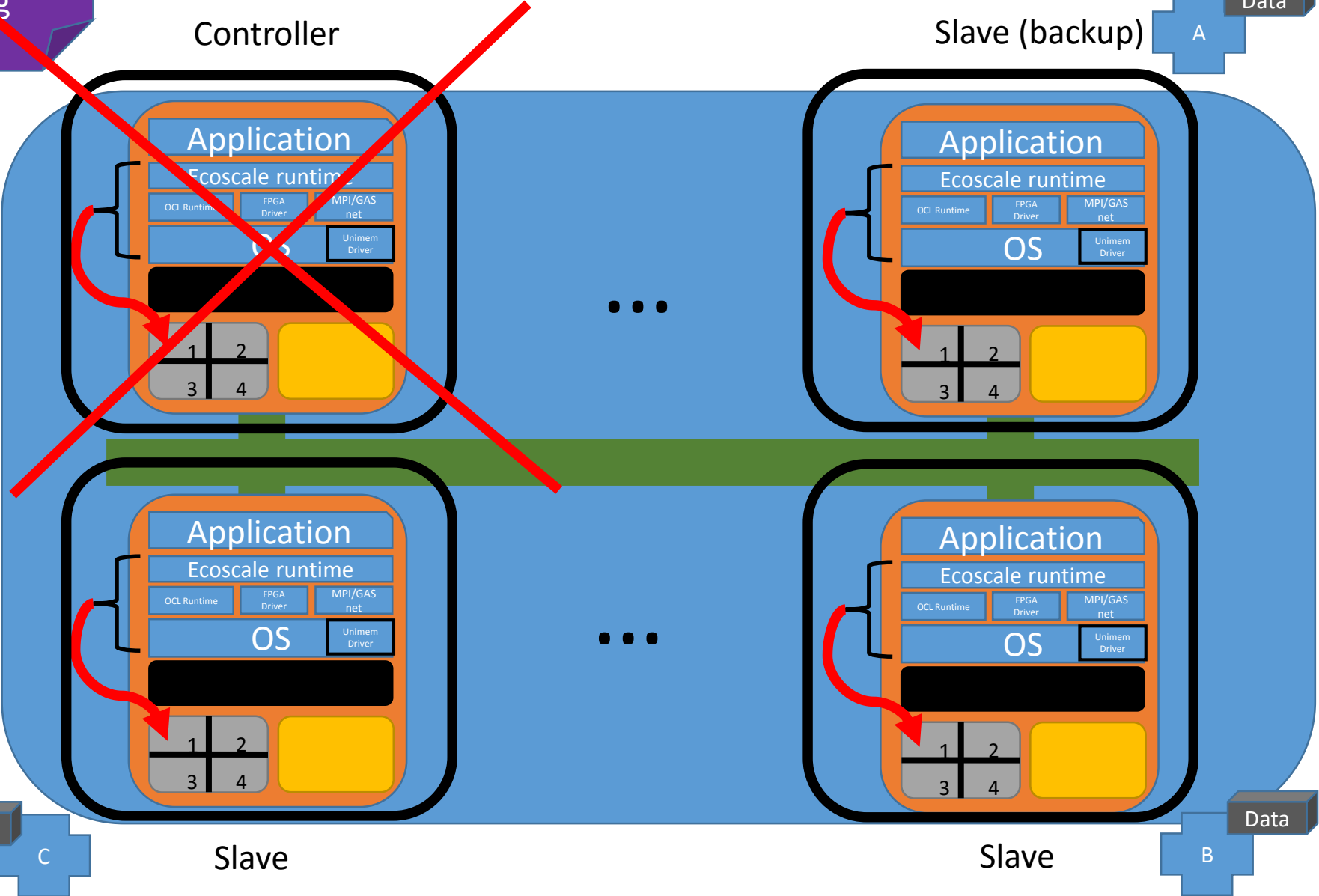
Worker Node

Unimem

CPU

FPGA

RAM



Accounting Log

Controller

Slave (backup)

Data

A

Data

C

Slave

Slave

B

Data

Compute Node

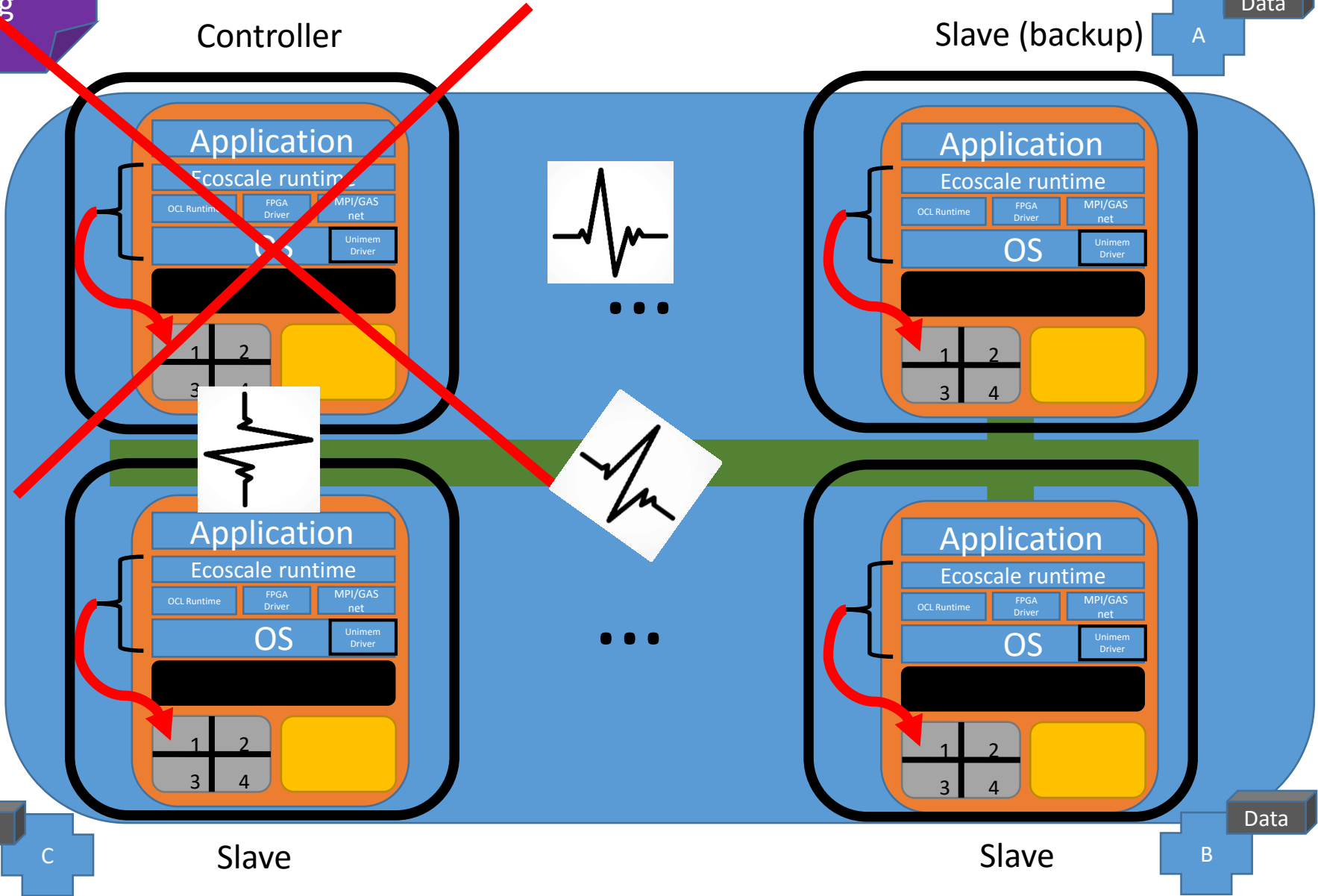
Worker Node

Unimem

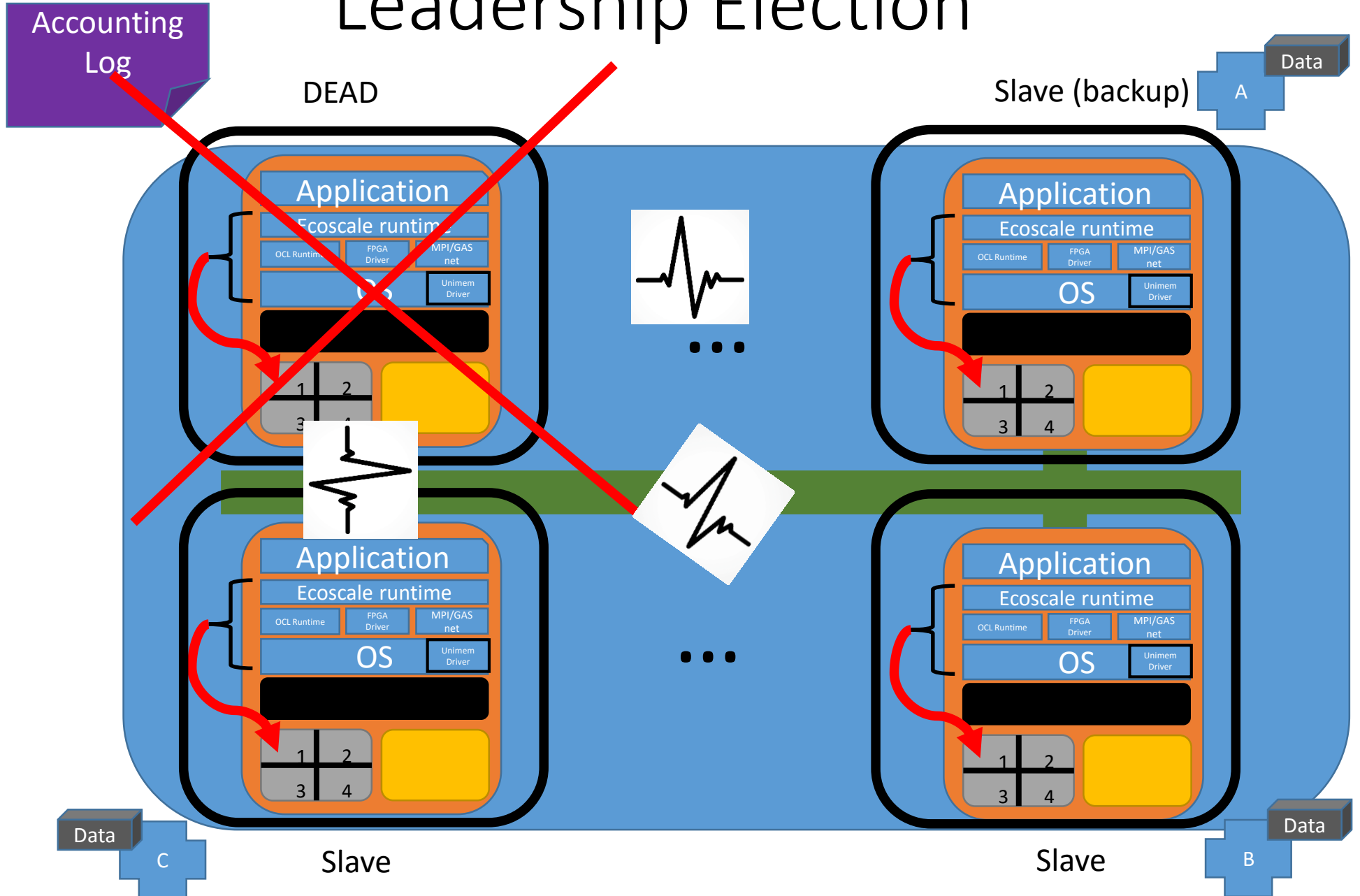
CPU

FPGA

RAM



Leadership Election



Compute Node

Worker Node

Unimem

CPU

FPGA

RAM

Accounting Log

DEAD

Controller

Data A

Compute Node

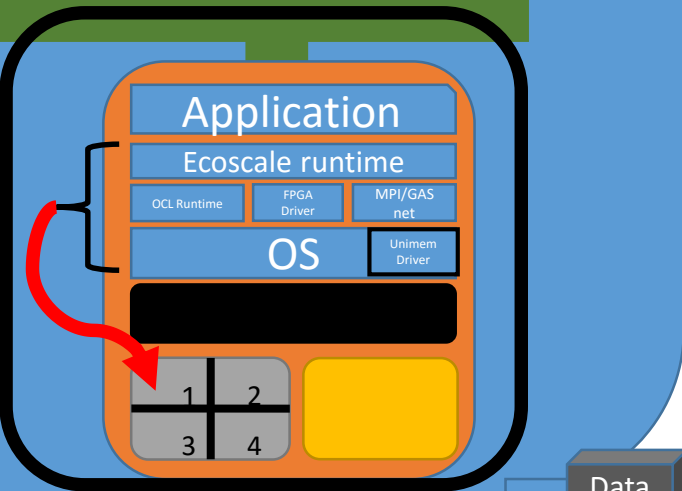
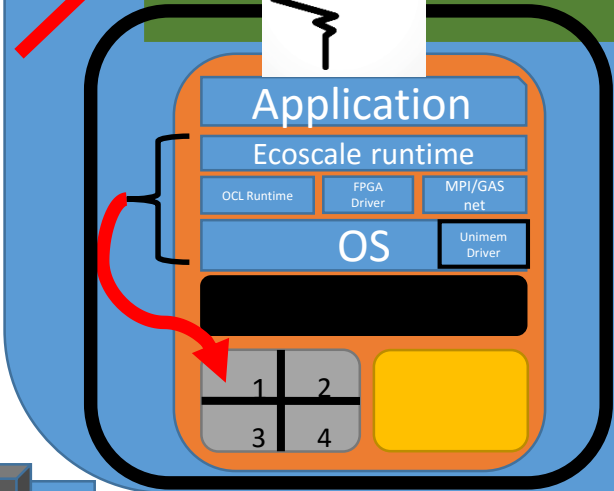
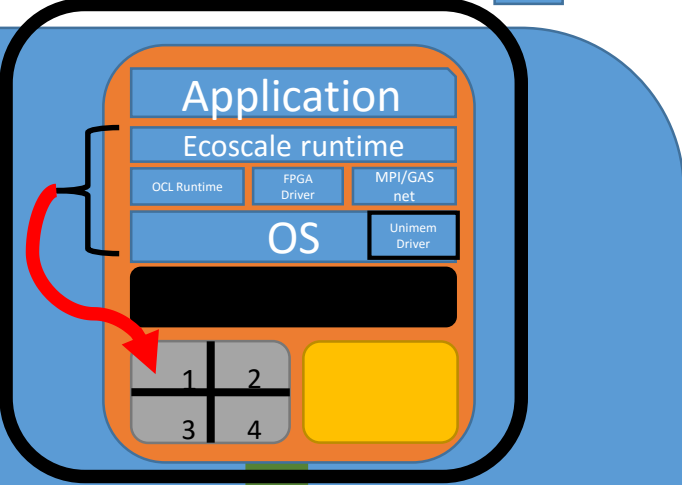
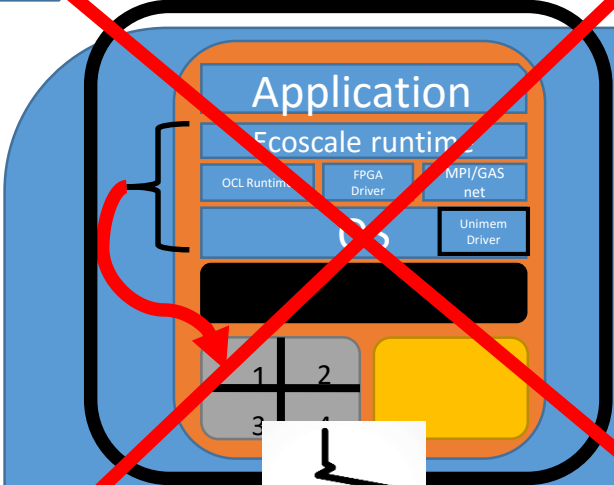
Worker Node

Unimem

CPU

FPGA

RAM



Data C

Slave (backup)

Slave

Data B

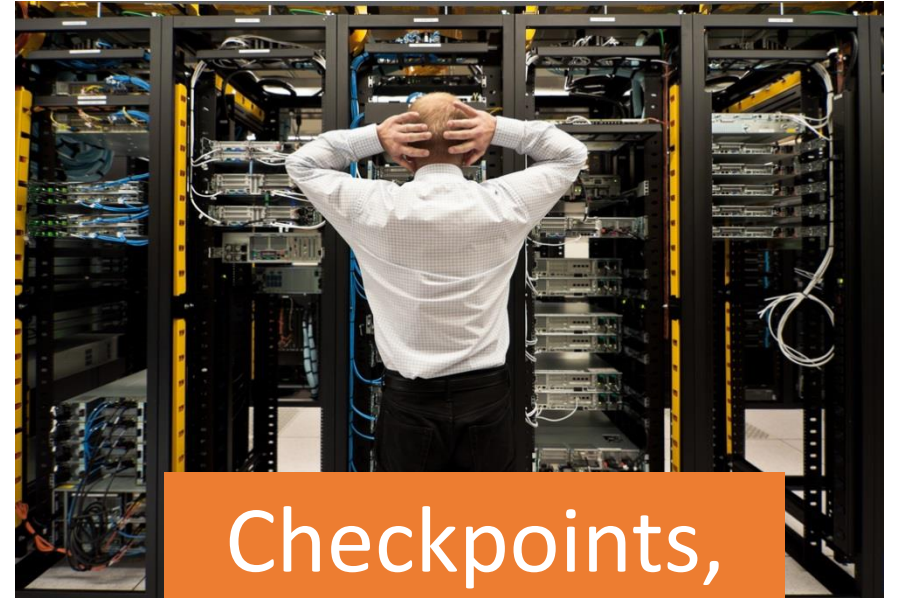
Exascale: Problems Solved?



Extended
OpenCL



FPGA



Checkpoints,
Heartbeats,
and internal
monitors

Ideas?

ありがとうございました!



質問はありますか



@jhebus



Paul-Harvey.org