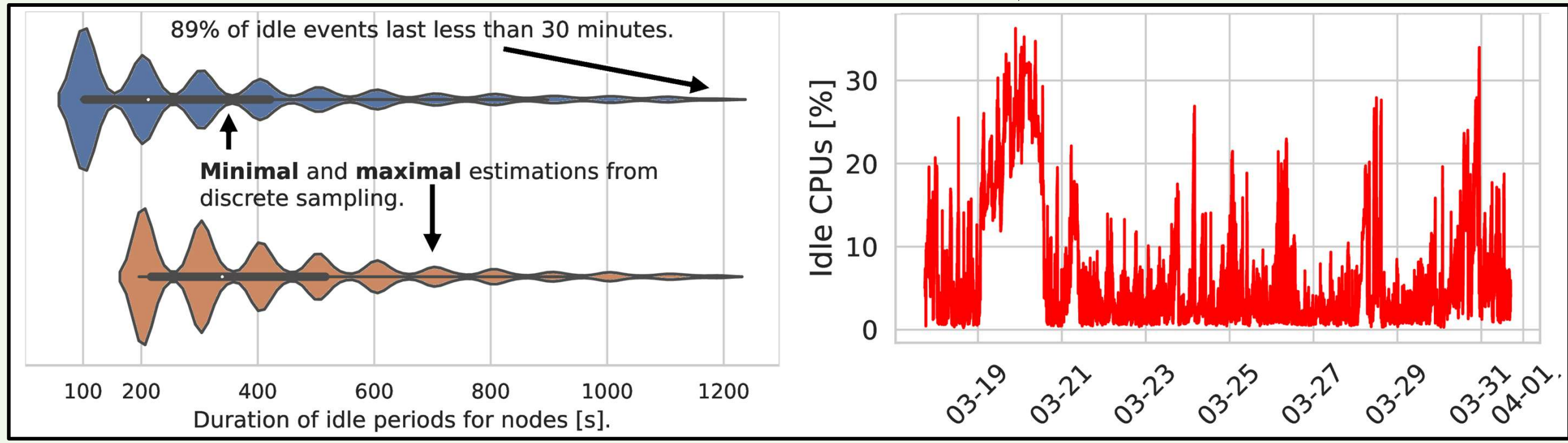


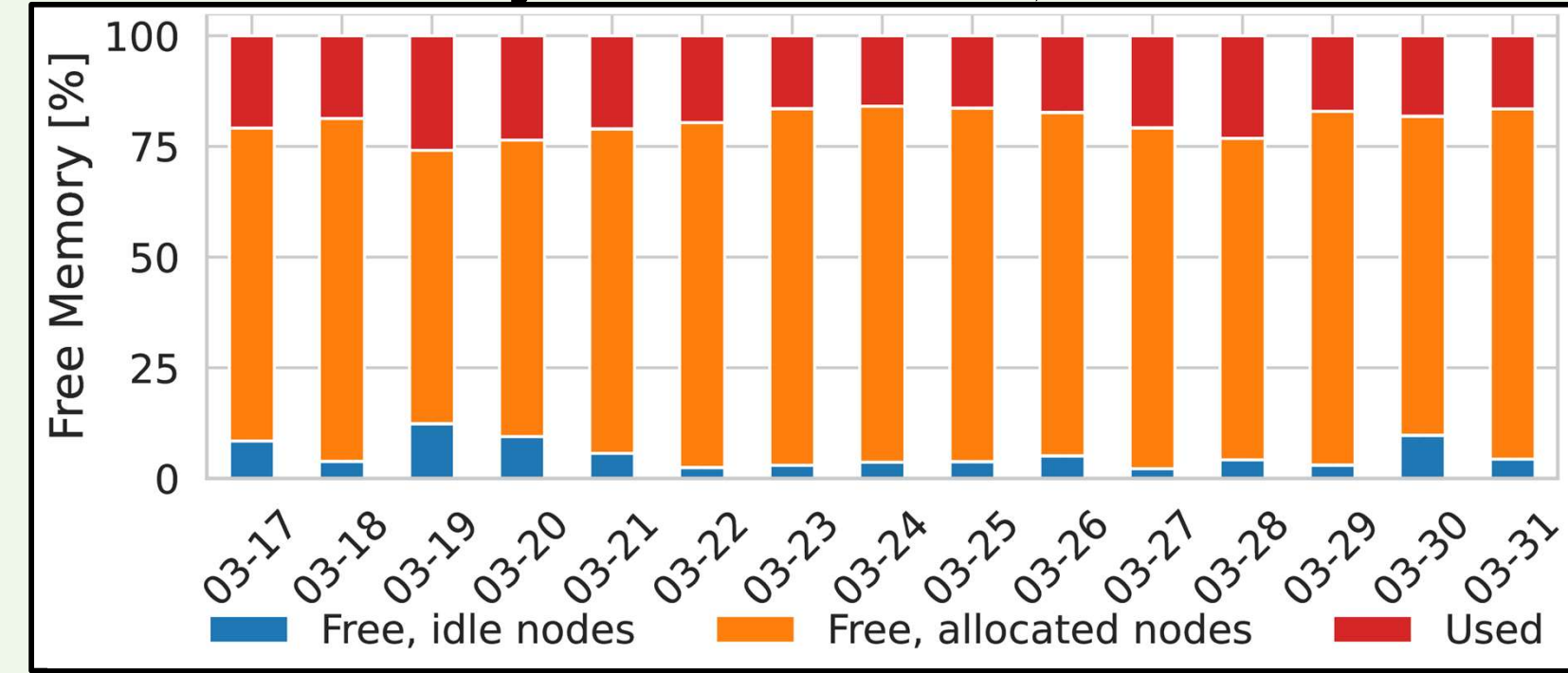
Challenge: Resource Underutilization in HPC Datacenters

Node and CPU Utilization, Piz Daint



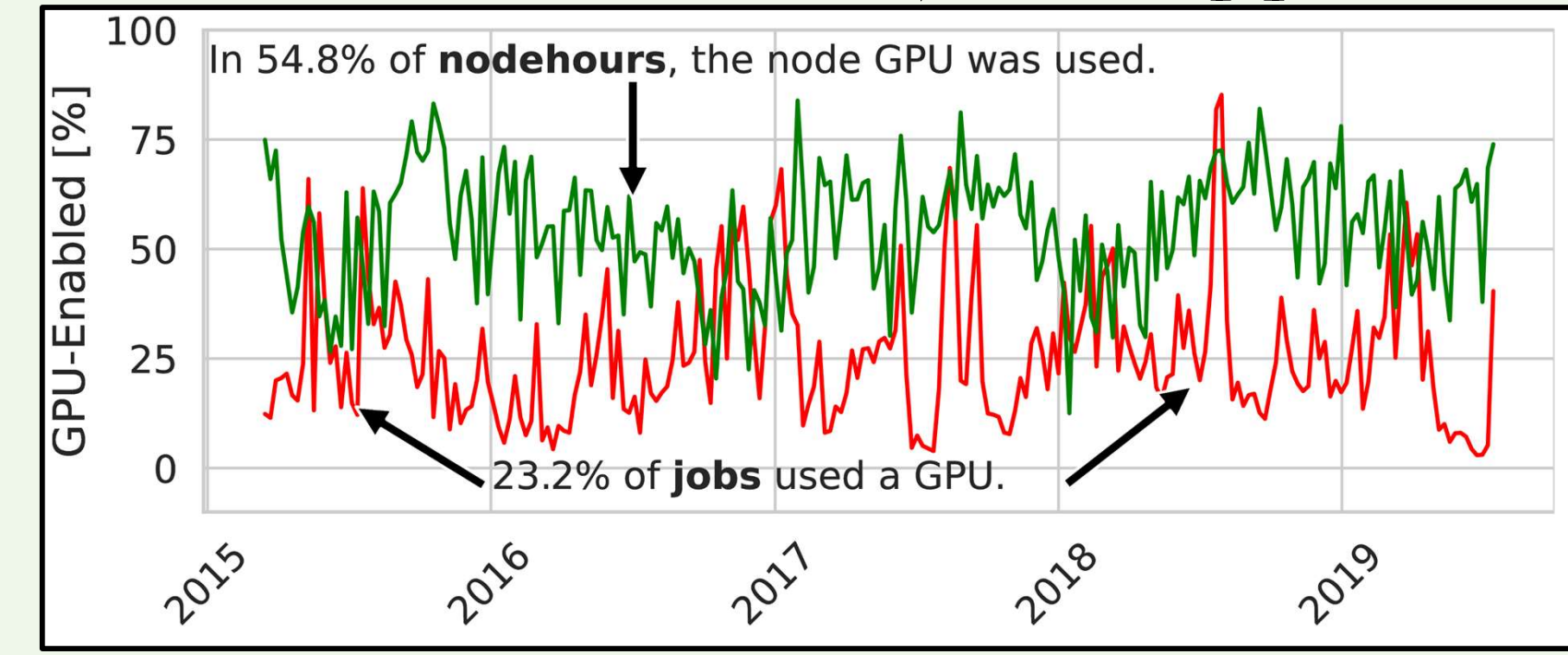
Nodes don't stay idle for an extended time – 70–80% are idle for less than 10 minutes. Long-running allocations cannot address these utilization gaps.

Memory Utilization, Piz Daint



Static allocations on homogeneous resources cannot improve memory utilization because these do not represent the heterogeneity of HPC workloads.

GPU Utilization, Titan [1]

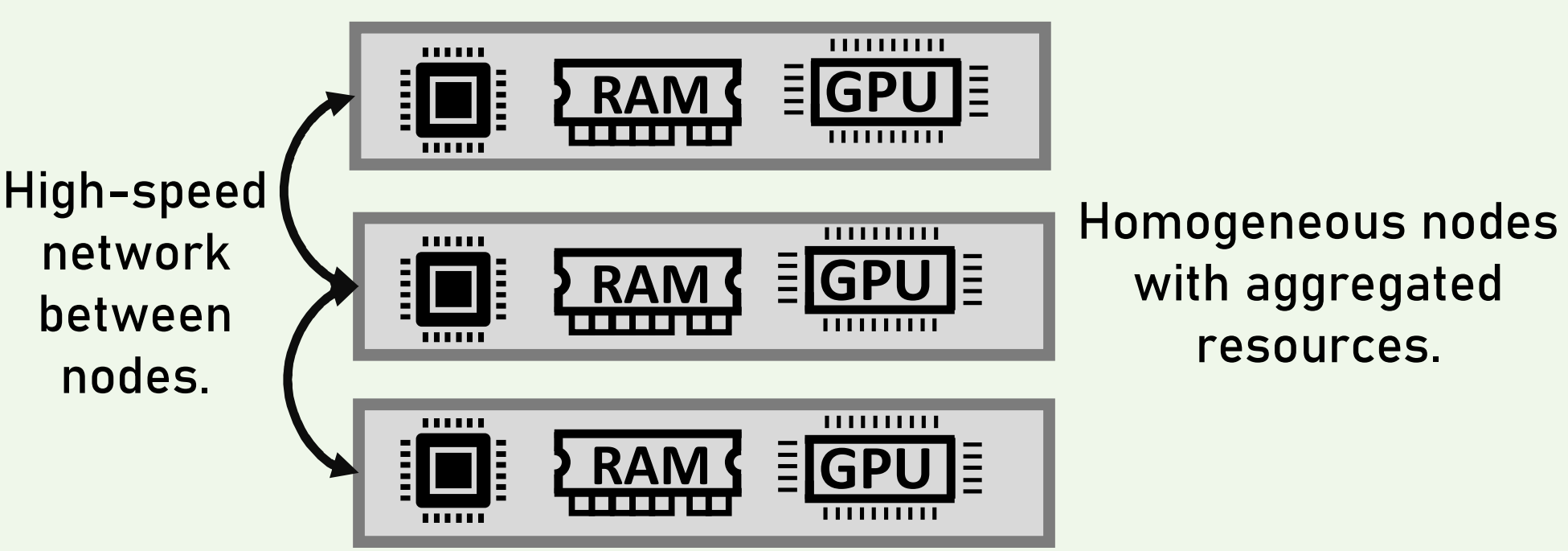


While HPC systems are getting more heterogeneous over time, GPU utilization by jobs is low, reinforcing the need to co-locate jobs.

To improve utilization of supercomputers, we need to enable sharing resources with fine-grained and short-term allocations.

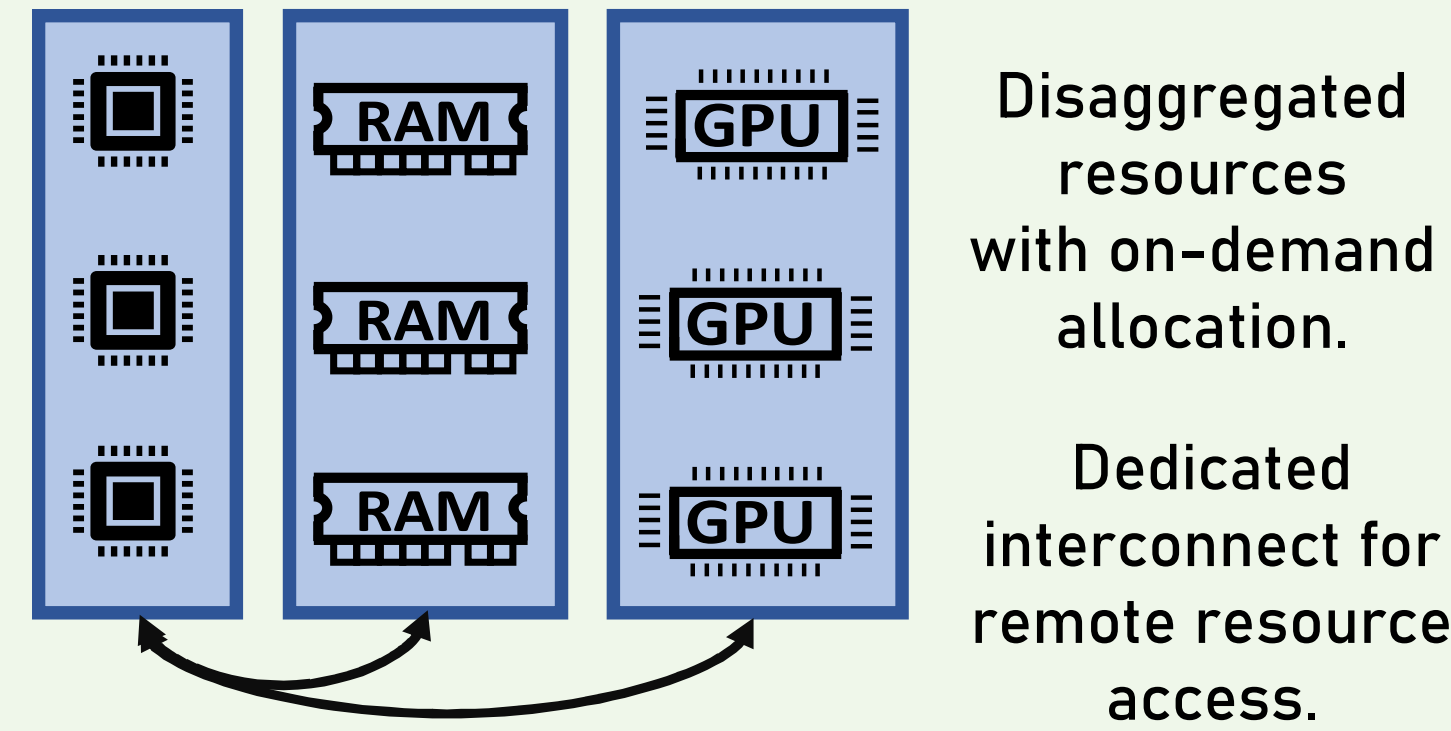
Solution: Software Resource Disaggregation with Serverless Functions

HPC Node – Tightly Coupled Hardware



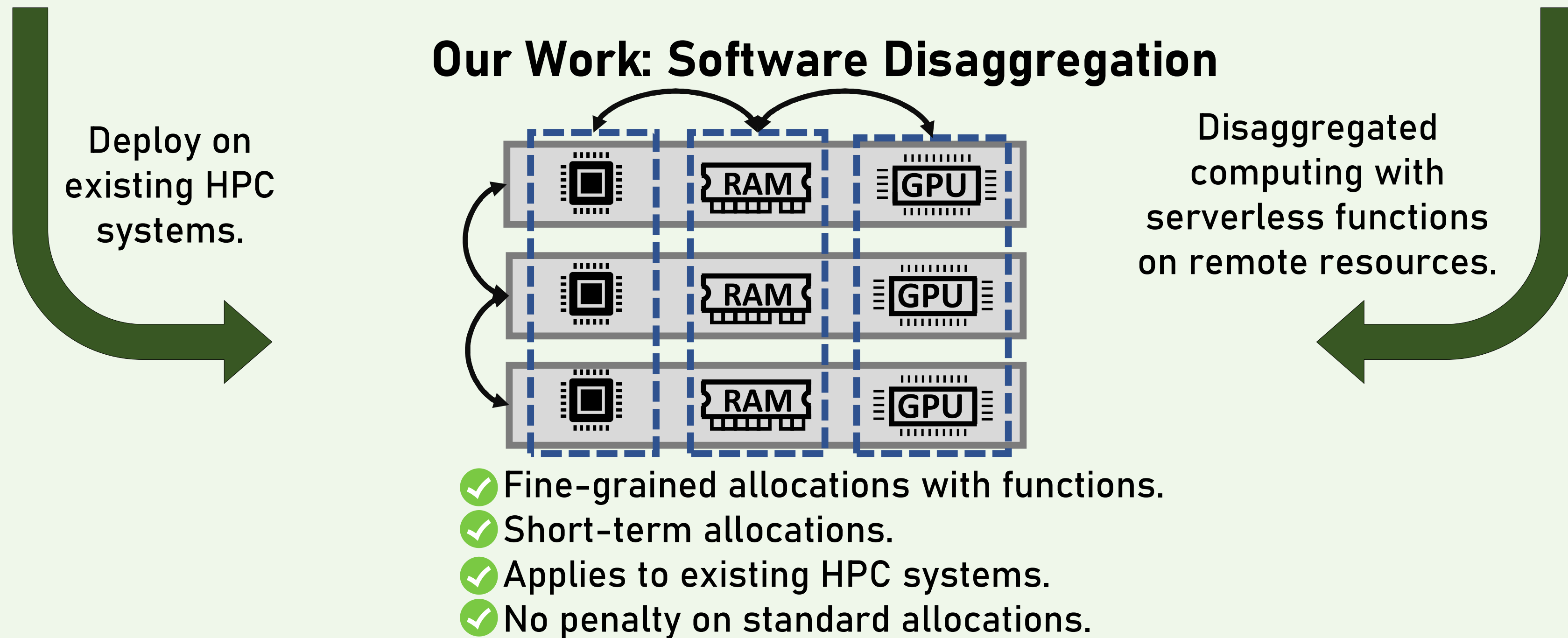
- ✓ No latency accessing resources.
- ✓ Allows for node sharing and job co-location...
- ✗ ... when resource consumption is compatible [2].
- ✗ Nodes are overprovisioned to support all jobs.
- ✗ No support for short allocations.

Hardware Disaggregated Data Center



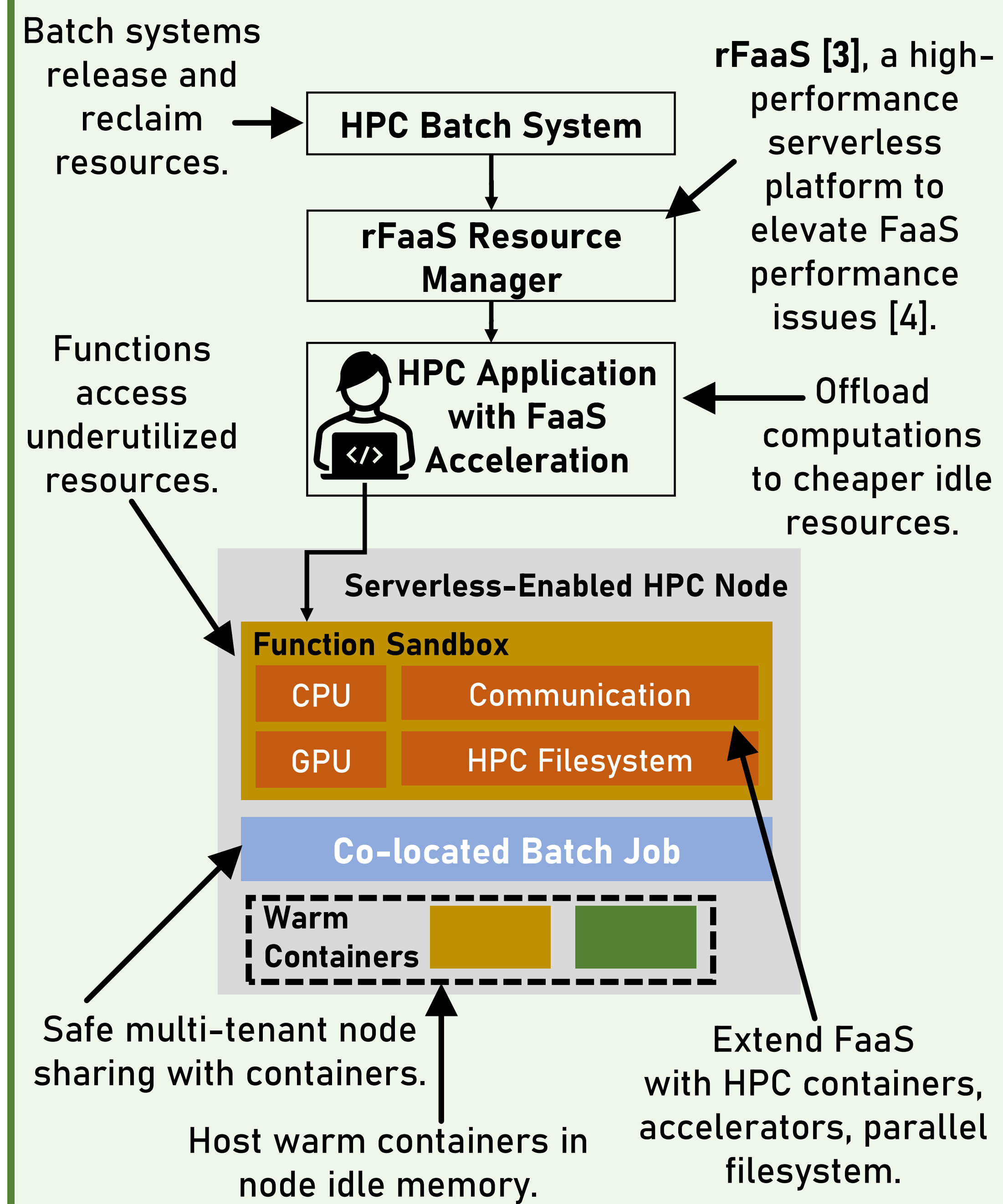
- ✓ Higher resource utilization.
- ✗ Requires new, dedicated hardware.
- ✗ Latency and bandwidth penalty.

Our Work: Software Disaggregation



- ✓ Fine-grained allocations with functions.
- ✓ Short-term allocations.
- ✓ Applies to existing HPC systems.
- ✓ No penalty on standard allocations.

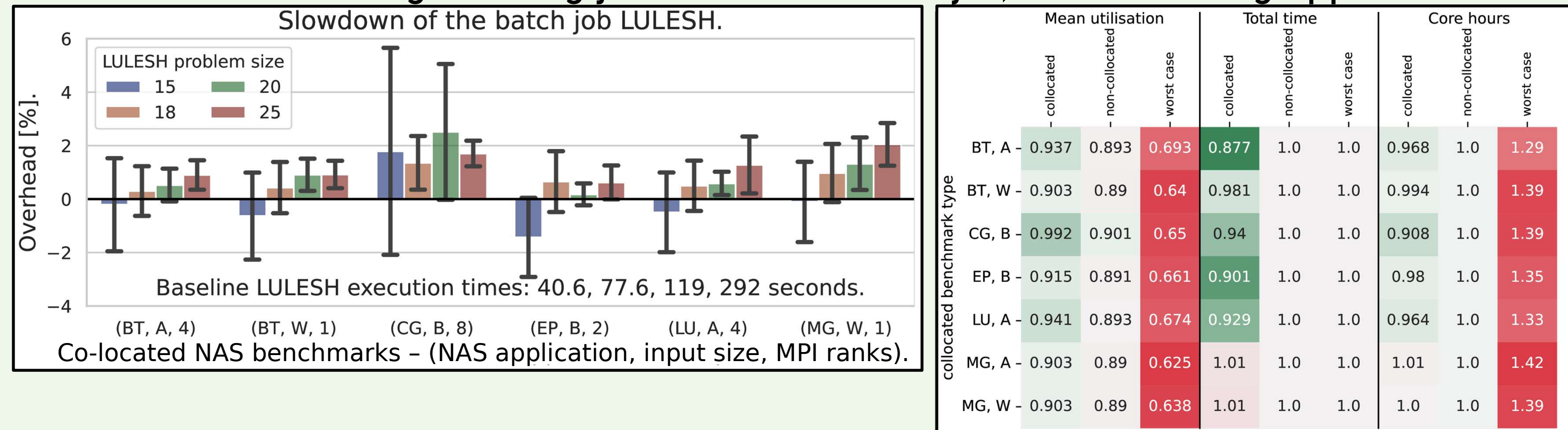
Bringing Serverless Disaggregation to HPC Systems with rFaaS



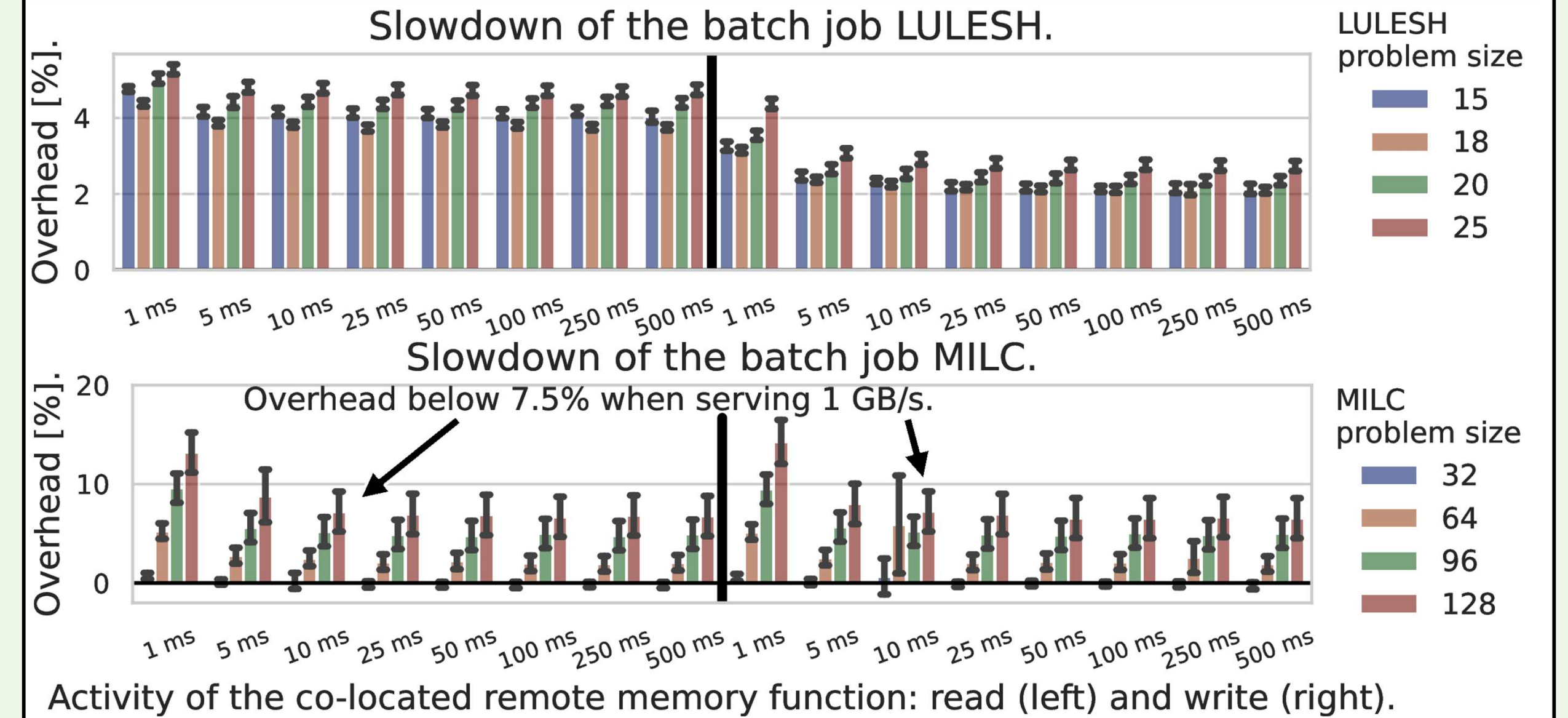
Evaluation

CPU Disaggregation: co-locating CPU workloads

Co-location of long-running job with function-style, short-running applications.

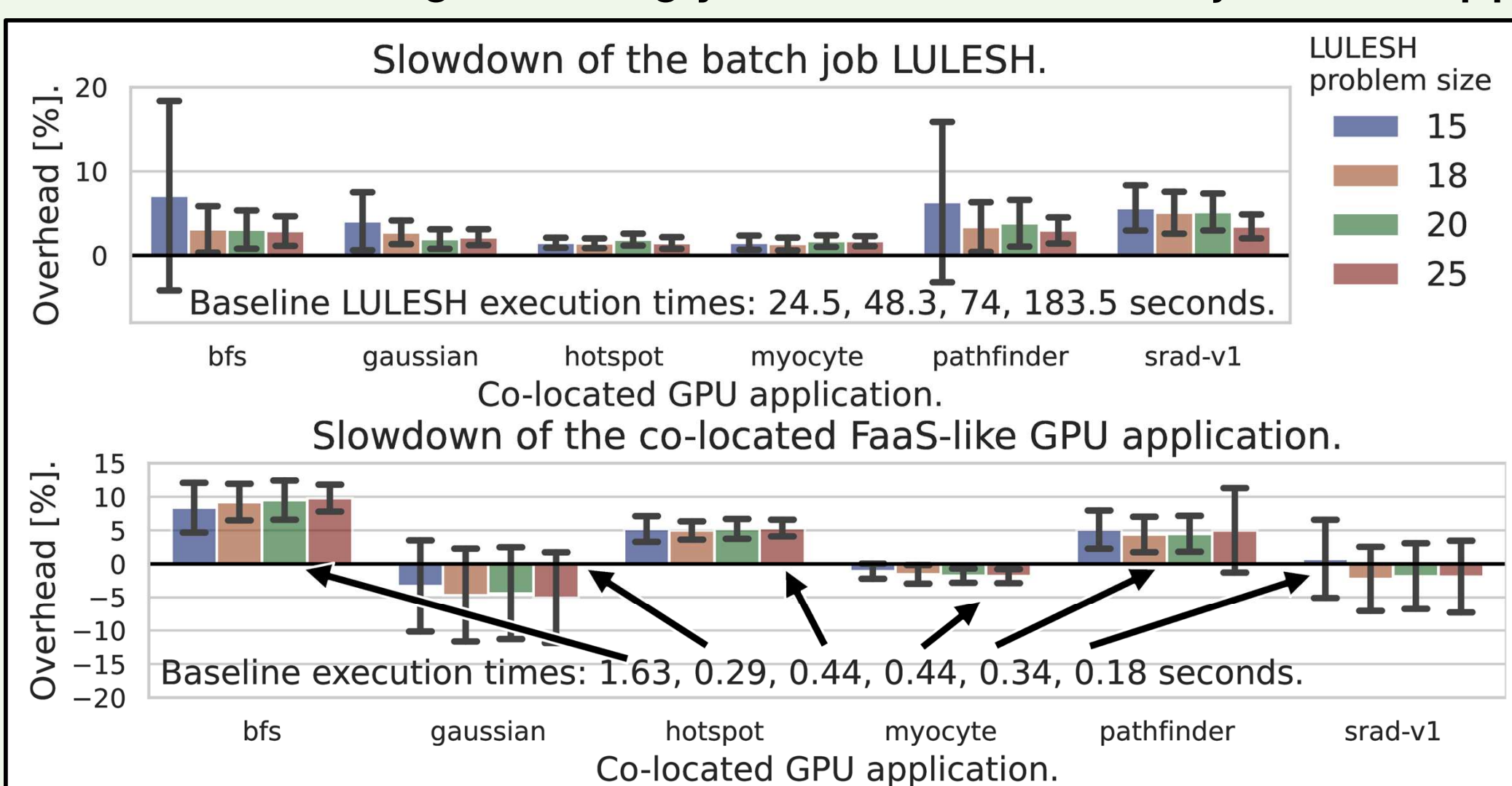


Memory disaggregation: co-locating RMA function



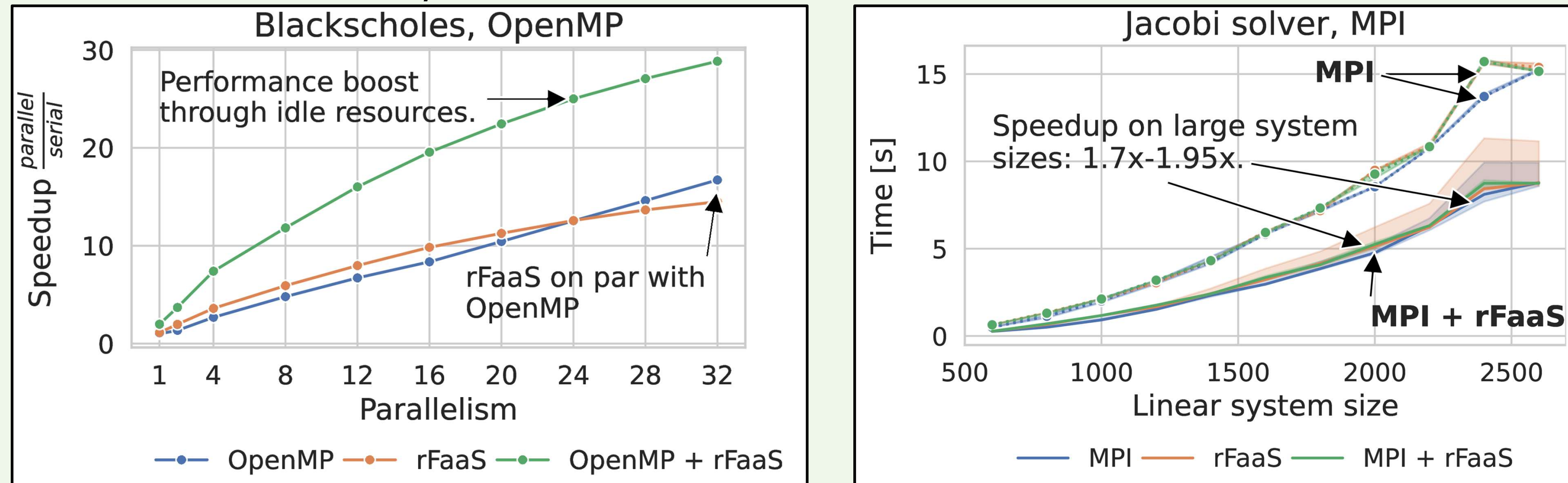
GPU disaggregation: co-locating GPU and CPU workloads

Co-location of long-running job with function-style GPU applications.



Integration of FaaS into HPC applications

Each thread/rank offloads half of their workload to an rFaaS function.



References

- [1] Wang F. et al., "Learning from Five-year Resource-Utilization Data of Titan System", IEEE CLUSTER 2019
- [2] Breslow A. et al., "The case for colocation of high performance computing workloads", Concurrency and Computation: Practice and Experience, 2013
- [3] Copik M. et al., "rFaaS: RDMA-Enabled FaaS Platform for Serverless High-Performance Computing", arXiv 2021
- [4] Copik M. et al., "SeBS: A Serverless Benchmark Suite for Function-as-a-Service Computing", ACM/IFIP Middleware 2021

