# Software Resource Disaggregation for HPC with Serverless Computing

**Marcin Copik,** Alexandru Calotoiu (advisor), Torsten Hoefler (advisor)

# Tracking Wasted Money in HPC

## Job Characteristics on Large-Scale Systems: Long-Term Analysis, Quantification, and Implications*

Tirthak Patel
Northeastern University

Zhengchun Liu, Raj Kettimuthu
Argonne National Laboratory

Paul Rich, William Allcock
Argonne National Laboratory

Devesh Tiwari
Northeastern University

SC, 2020

## FINAL REPORT
## WORKLOAD ANALYSIS OF BLUE WATERS
### (ACI 1650758)

Matthew D. Jones, Joseph P. White, Martins Innus, Robert L. DeLeon, Nikolay Simakov, Jeffrey T. Palmer, Steven M. Gallo, and Thomas R. Furlani (furlani@buffalo.edu), Center for Computational Research, University at Buffalo, SUNY

Michael Showerman, Robert Brunner, Andry Kot, Gregory Bauer, Brett Bode, Jeremy Enos, and William Kramer (wtkramer@illinois.edu), National Center for Supercomputing Applications (NCSA), University of Illinois at Urbana Champaign

arXiv, 2017

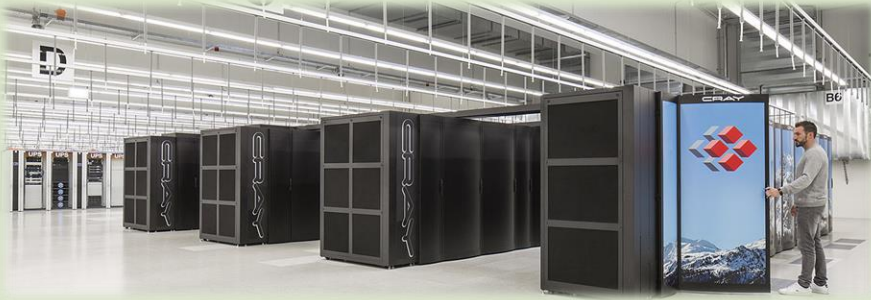## Comprehensive Workload Analysis and Modeling of a Petascale Supercomputer

Haihang You[1] and Hao Zhang[2]

[1] National Institute for Computational Sciences,
Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA
[2] Department of Electrical Engineering and Computer Science,
University of Tennessee, Knoxville, TN 37996, USA
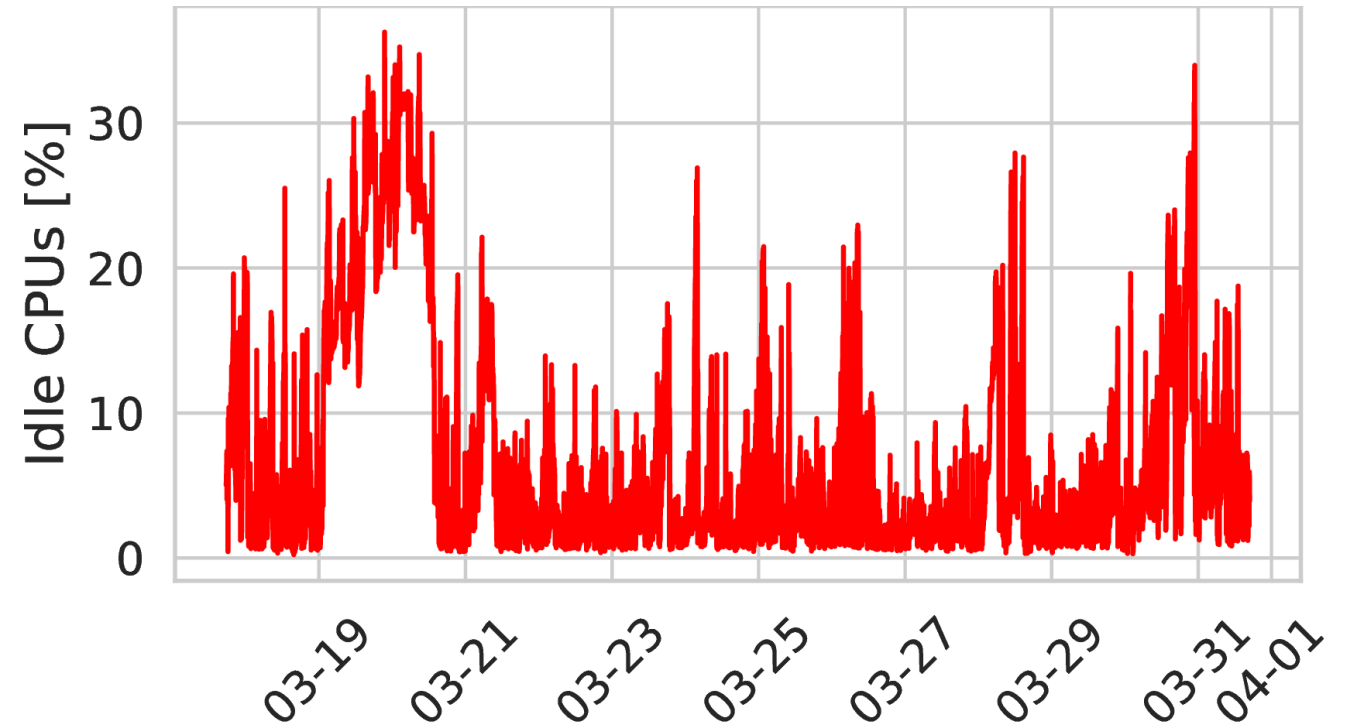{hyou,haozhang}@utk.edu

JSSPP, 2012

2

# HPC System Utilization - CPU



**Piz Daint,** April 2022.
- XC50 nodes – CPU + GPU, 64 GB memory.
- XC40 nodes – CPU, 64/128 GB memory.

Query SLURM info every two minutes.
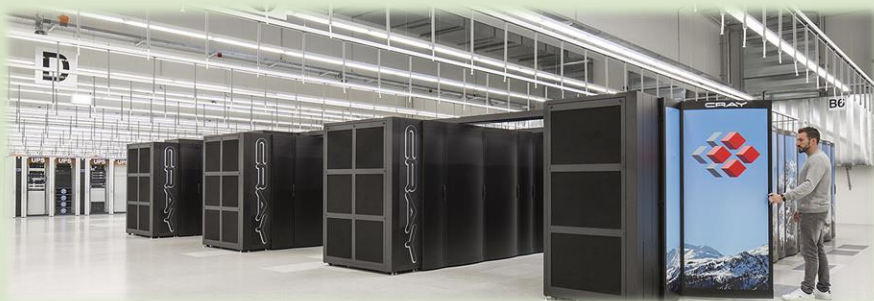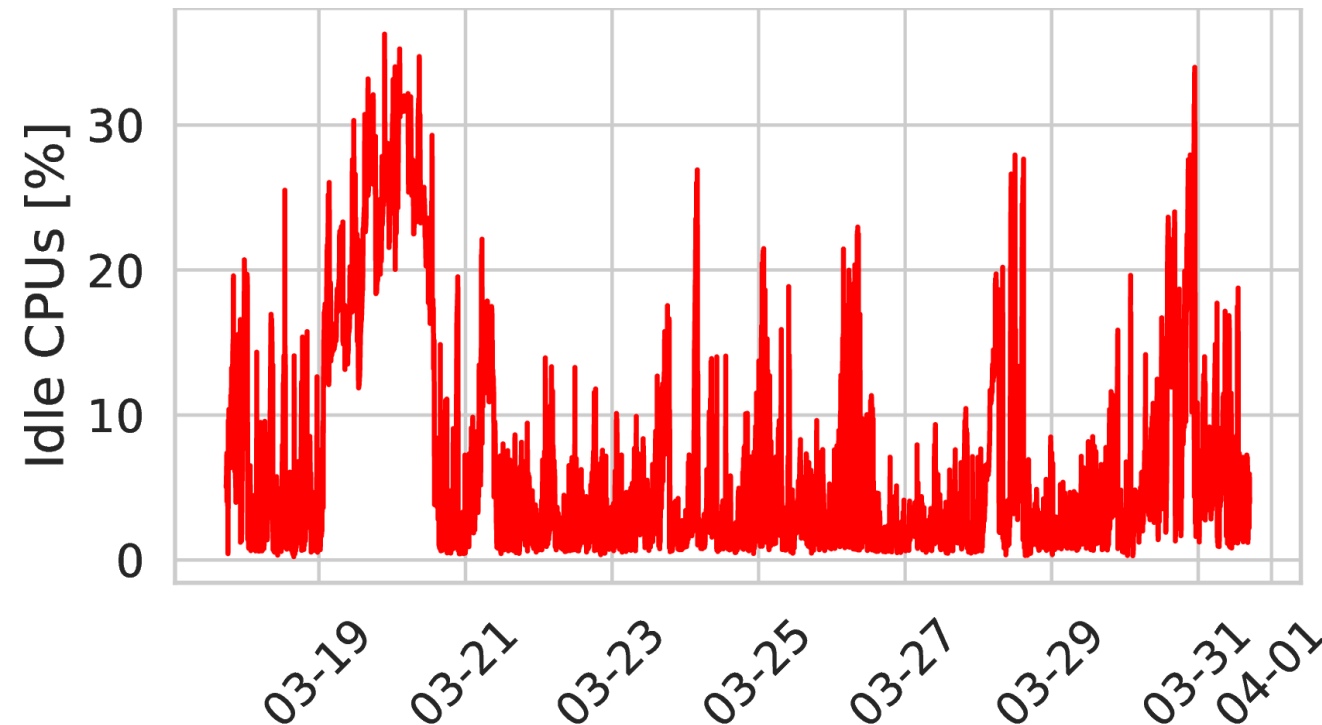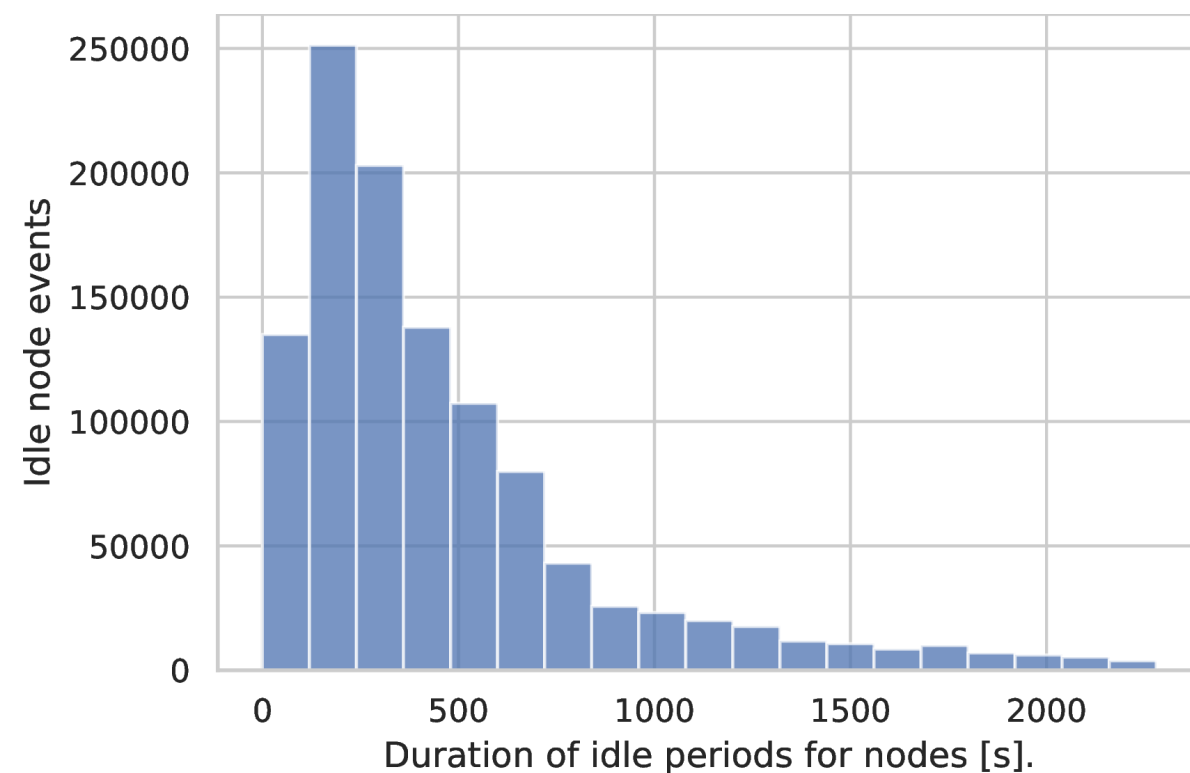
# HPC System Utilization - CPU

**Piz Daint,** April 2022.
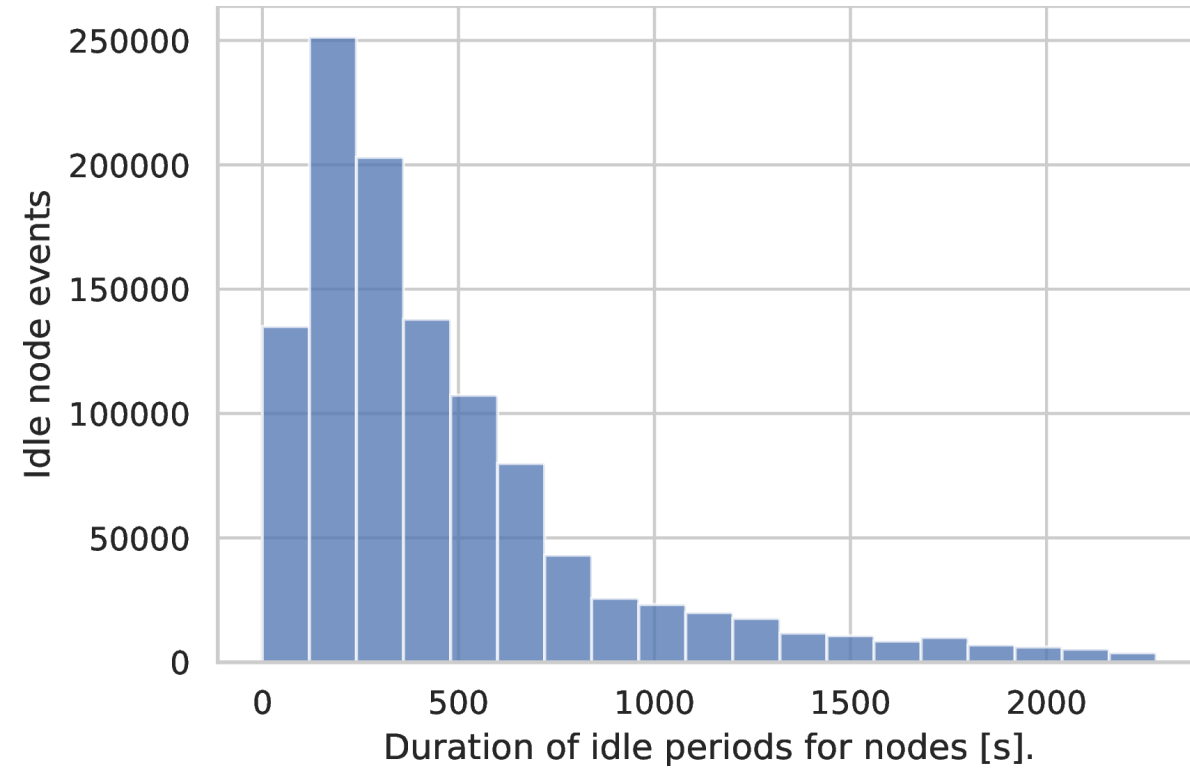- XC50 nodes – CPU + GPU, 64 GB memory.
- XC40 nodes – CPU, 64/128 GB memory.

Query SLURM info every two minutes.

# HPC System Utilization - CPU



**Piz Daint,** April 2022.
- XC50 nodes – CPU + GPU, 64 GB memory.
- XC40 nodes – CPU, 64/128 GB memory.

Query SLURM info every two minutes.



# How long do nodes stay idle?

# HPC System Utilization - CPU

## Minimum

## Maximum

# HPC System Utilization - CPU

# HPC System Utilization - CPU

## Minimum



## Maximum



80% and 70% of idle node events last less than 10 minutes.
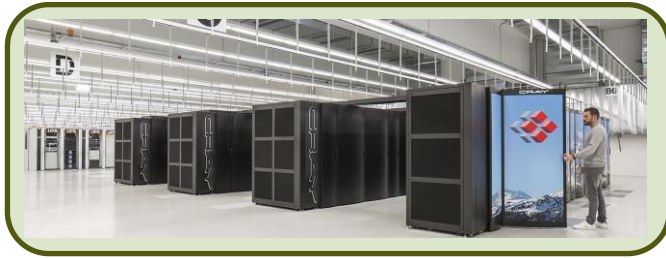
# HPC System Utilization - CPU



## Minimum

## Maximum

**Short-term resource availability requires short-term allocations.**

80% and 70% of idle node events last less than 10 minutes.

# HPC System Utilization - Memory





Free Memory [%]

100 · 75 · 50 · 25 · 0

03-17 · 03-18 · 03-19 · 03-20 · 03-21 · 03-22 · 03-23 · 03-24 · 03-25 · 03-26 · 03-27 · 03-28 · 03-29 · 03-30 · 03-31

Free, idle nodes    Free, allocated nodes    Used

# HPC System Utilization - Memory

# HPC System Utilization - Memory

## A Case For Intra-rack Resource Disaggregation in HPC

GEORGE MICHELOGIANNAKIS, Lawrence Berkeley National Laboratory, USA
BENJAMIN KLENK, NVIDIA, USA
BRANDON COOK, Lawrence Berkeley National Laboratory, USA
MIN YEE TEH and MADELEINE GLICK, Columbia University, USA
LARRY DENNISON, NVIDIA, USA
KEREN BERGMAN, Columbia University, USA
JOHN SHALF, Lawrence Berkeley National Laboratory, USA

TACO, 2022

## Quantifying Memory Underutilization in HPC Systems and Using it to Improve Performance via Architecture Support

Gagandeep Panwar[*]
Virginia Tech
Blacksburg, USA
gpanwar@vt.edu

Da Zhang[*]
Virginia Tech
Blacksburg, USA
daz3@vt.edu

Yihan Pang[*]
Virginia Tech
Blacksburg, USA
pyihan1@vt.edu

Mai Dahshan
Virginia Tech
Blacksburg, USA
mdahshan@vt.edu

Nathan DeBardeleben
Los Alamos National Laboratory
Los Alamos, USA
ndebard@lanl.gov

Binoy Ravindran
Virginia Tech
Blacksburg, USA
binoy@vt.edu

Xun Jian
Virginia Tech
Blacksburg, USA
xunj@vt.edu

MICRO, 2019

## FINAL REPORT
## WORKLOAD ANALYSIS OF BLUE WATERS
### (ACI 1650758)

Matthew D. Jones, Joseph P. White, Martins Innus, Robert L. DeLeon, Nikolay Simakov, Jeffrey T. Palmer, Steven M. Gallo, and Thomas R. Furlani (furlani@buffalo.edu), Center for Computational Research, University at Buffalo, SUNY

Michael Showerman, Robert Brunner, Andry Kot, Gregory Bauer, Brett Bode, Jeremy Enos, and William Kramer (wtkramer@illinois.edu), National Center for Supercomputing Applications (NCSA), University of Illinois at Urbana Champaign

arXiv, 2017

## A Holistic View of Memory Utilization on HPC Systems: Current and Future Trends

Ivy B. Peng[*]
peng8@llnl.gov
Lawrence Livermore National Laboratory
USA

Ian Karlin
karlin1@llnl.gov
Lawrence Livermore National Laboratory
USA

Maya B. Gokhale
gokhale2@llnl.gov
Lawrence Livermore National Laboratory
USA

Kathleen Shoga
Shoga1@llnl.gov
Lawrence Livermore National Laboratory
USA

Matthew Legendre
legendre1@llnl.gov
Lawrence Livermore National Laboratory
USA

Todd Gamblin
gamblin2@llnl.gov
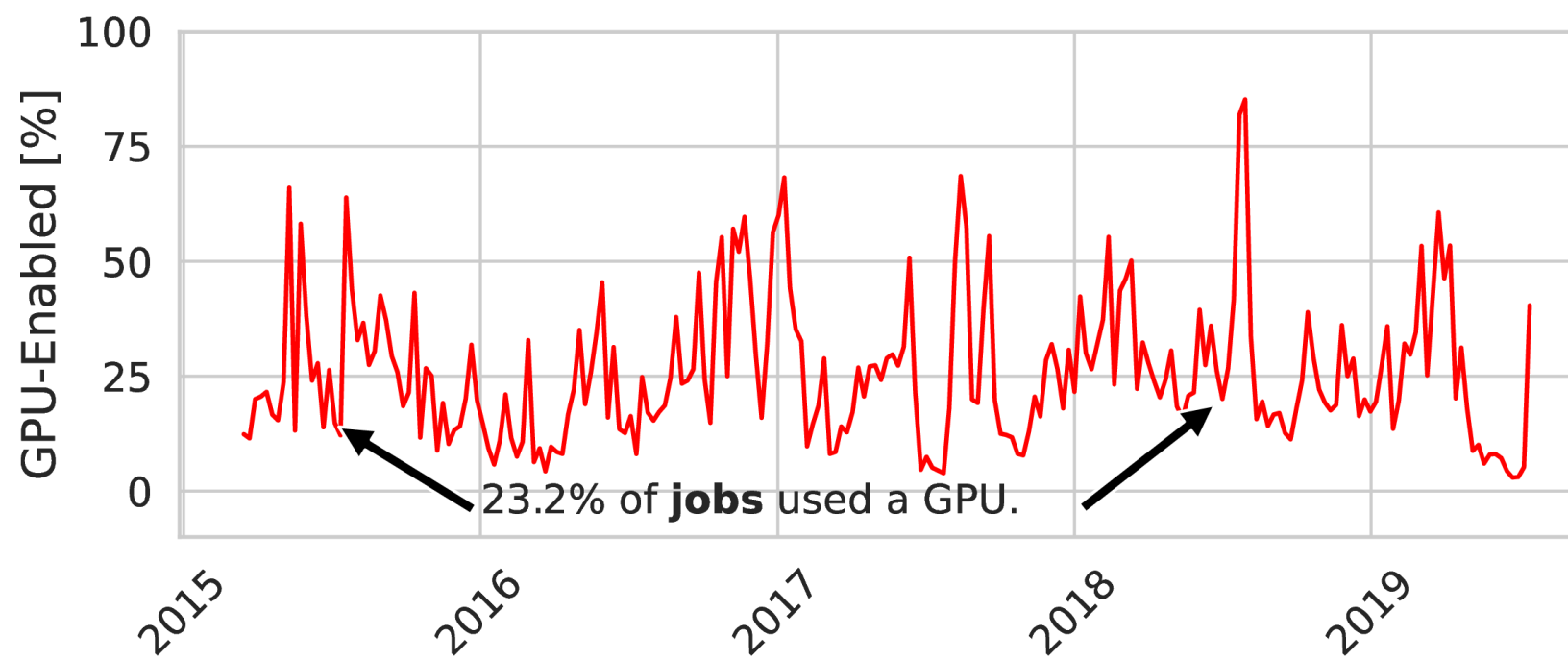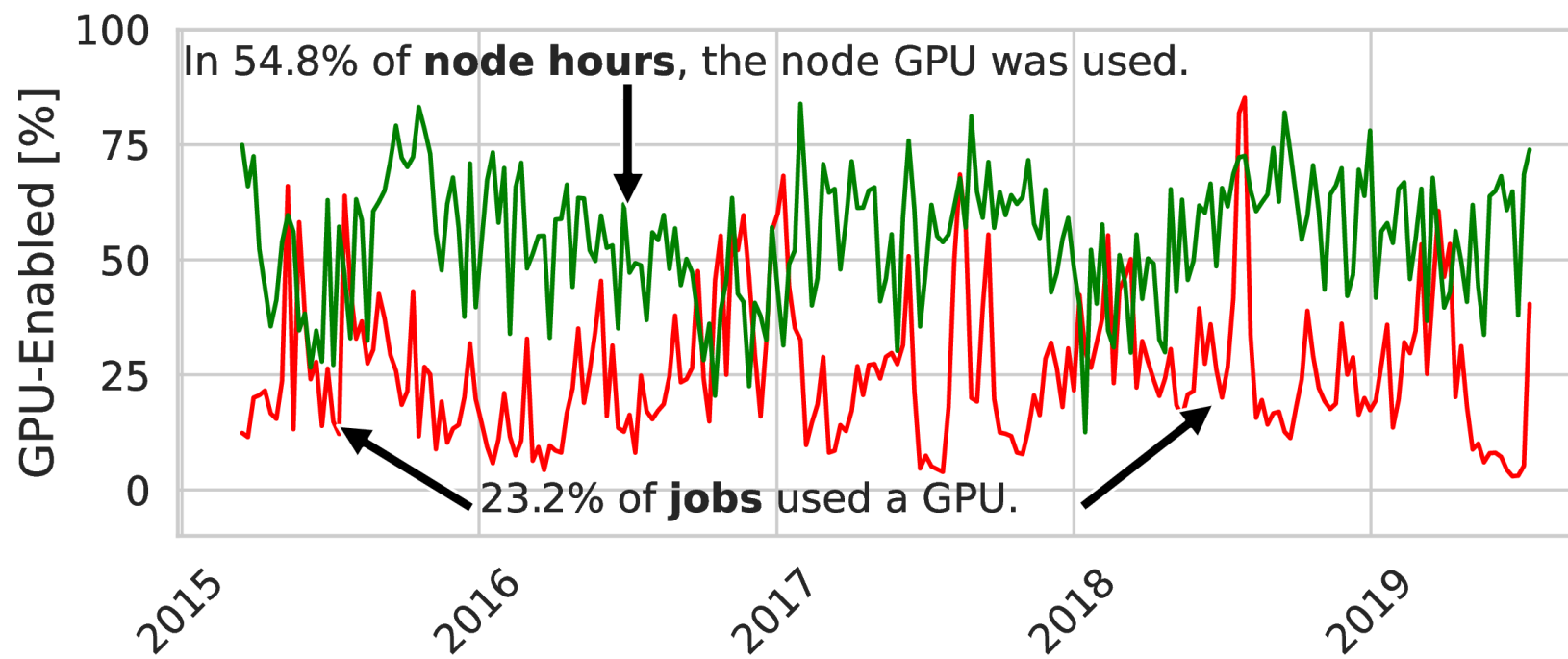Lawrence Livermore National Laboratory
USA

MEMSYS, 2021

# HPC System Utilization - GPU

Learning from Five-year Resource-Utilization Data
of Titan System

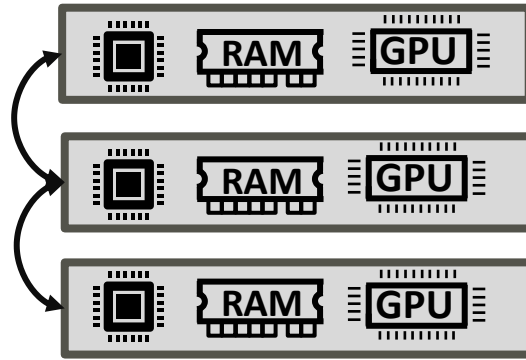Feiyi Wang[*], Sarp Oral[†], Satyabrata Sen[‡] and Neena Imam[§]
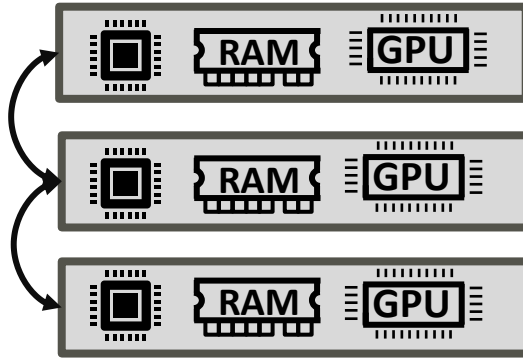
Oak Ridge National Laboratory

CLUSTER, 2019

# HPC System Utilization - GPU

23.2% of **jobs** used a GPU.

# HPC System Utilization - GPU

Learning from Five-year Resource-Utilization Data
of Titan System

Feiyi Wang[*], Sarp Oral[†], Satyabrata Sen [‡] and Neena Imam[§]

Oak Ridge National Laboratory

CLUSTER, 2019



In 54.8% of **node hours**, the node GPU was used.

23.2% of **jobs** used a GPU.

# Software Solution

**Standard HPC Node**



✅ High performance

❌ Inflexible architecture

# Software Solution

**Standard HPC Node**



**Hardware Disaggregation**



✅ High performance

❌ Inflexible architecture

✅ High efficiency

❌ Cost, performance penalty

# Software Solution



**Standard HPC Node**

✅ High performance

❌ Inflexible architecture

**Existing Coupled
Hardware Systems**

**Hardware Disaggregation**

✅ High efficiency

❌ Cost, performance penalty

# Software Solution

**Standard HPC Node**



**Hardware Disaggregation**



✅ High performance

❌ Inflexible architecture

✅ High efficiency

❌ Cost, performance penalty

**Existing Coupled Hardware Systems**

**Software Abstraction for Disaggregation**

We propose a software disaggregation approach to share node resources

We propose a **software disaggregation** approach to **share node resources** between **coarse-grained, long-running, and static batch jobs**

We propose a **software disaggregation** approach to **share node resources** between
**coarse-grained, long-running, and static batch jobs**
and
**fine-grained, short-term, and dynamically allocated serverless functions.**

# Serverless as an Answer

# Serverless as an Answer

# Serverless as an Answer

## Hardware Abstraction

# Serverless as an Answer

**Hardware Abstraction**



Control Plane

RAM  GPU

# Serverless as an Answer

## Hardware Abstraction



## Software Abstraction

# Serverless as an Answer

## Hardware Abstraction



## Software Abstraction

# Serverless as an Answer

## Hardware Abstraction



**Pay-as-you-go billing**

## Software Abstraction



**Granular computing**

# Serverless Disaggregation

## Batch jobs

# Serverless Disaggregation
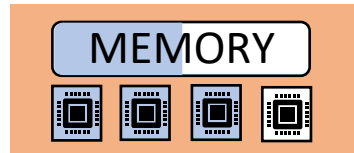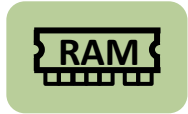
## Batch jobs

## Batch jobs + serverless functions

# Serverless Disaggregation

# Serverless Disaggregation
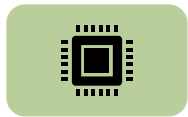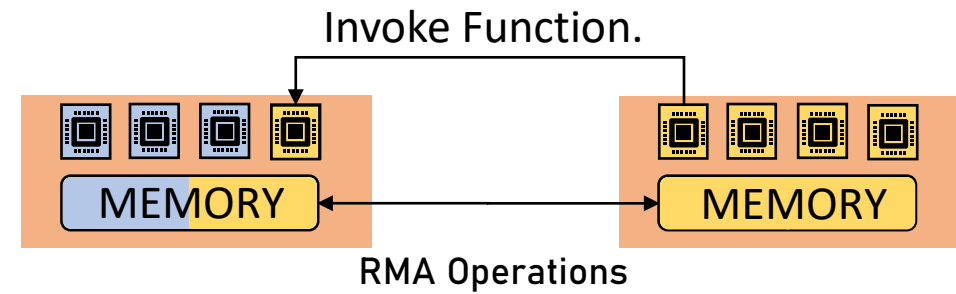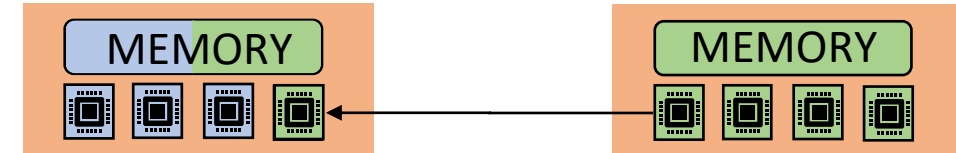
**Batch jobs**

**Batch jobs + serverless functions**

# Serverless Disaggregation

**Batch jobs**

**Batch jobs + serverless functions**

# Serverless Disaggregation

**Batch jobs**

**Batch jobs + serverless functions**
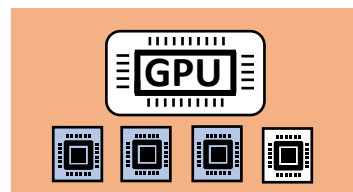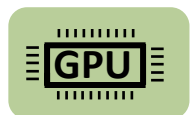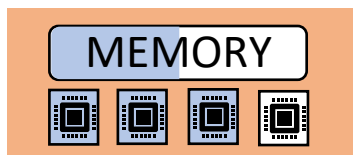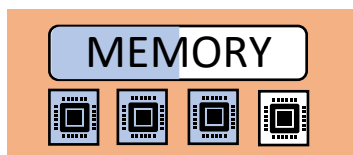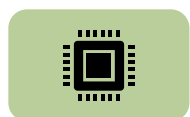


Invoke Function.

# Serverless Disaggregation

**Batch jobs**                    **Batch jobs + serverless functions**
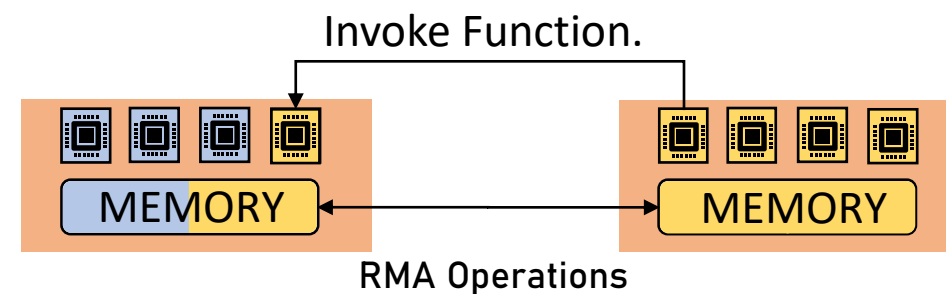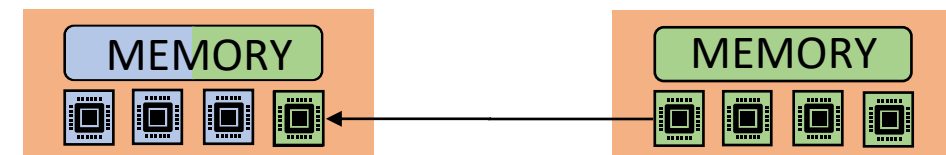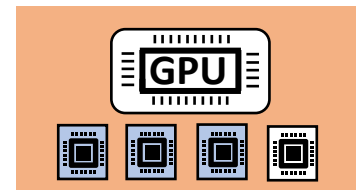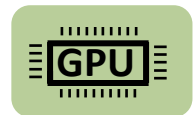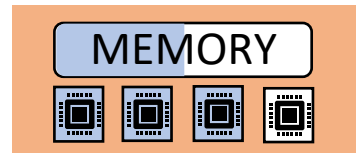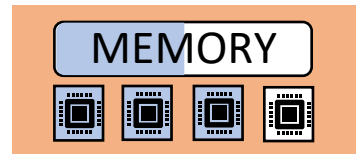


Invoke Function.
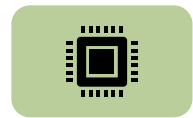
RMA Operations

# Serverless Disaggregation



Batch jobs

Batch jobs + serverless functions

Invoke Function.

RMA Operations

# Serverless Disaggregation

# Serverless Disaggregation

# Serverless Disaggregation

# Serverless Disaggregation

# Evaluation



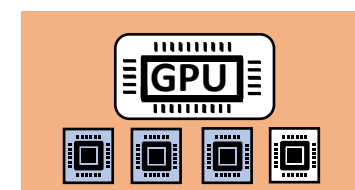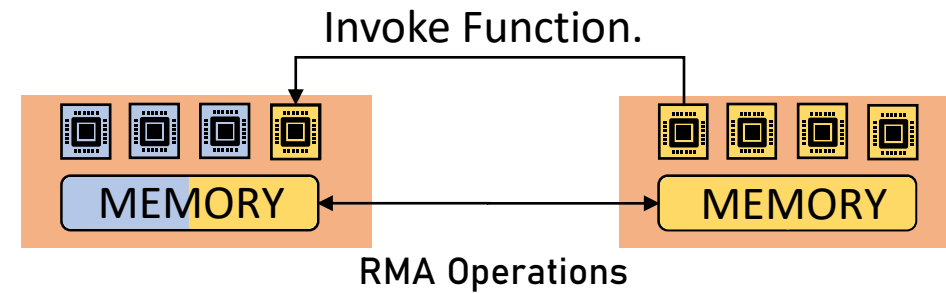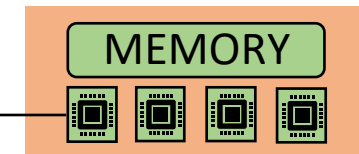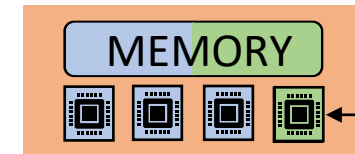**XC50** nodes - 12 CPU cores, GPU, 64 GB memory.

**XC40** nodes - 36 CPU cores, 64/128 GB memory.

**Cray Aries** interconnect.

36 CPU cores, 377 GB memory.
Ethernet with RoCEv2 support.

# #1 CPU Sharing



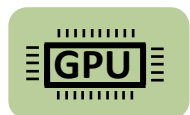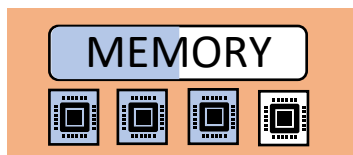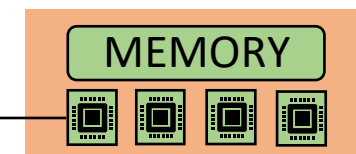| | Mean utilisation | | |
|---|---|---|---|
| collocated benchmark type | Disaggregation | Ideal Non-sharing | Realistic |
| BT, A | 0.937 | 0.893 | 0.693 |
| BT, W | 0.903 | 0.89 | 0.64 |
| CG, B | 0.992 | 0.901 | 0.65 |
| EP, B | 0.915 | 0.891 | 0.661 |
| LU, A | 0.941 | 0.893 | 0.674 |
| MG, A | 0.903 | 0.89 | 0.625 |
| MG, W | 0.903 | 0.89 | 0.638 |

**LULESH**
64 ranks, 2 nodes
32 out of 36 cores allocated.

**NAS**
1 – 4 ranks
Distributed across nodes.

# #1 CPU Sharing

|  | Mean utilisation | | | Total time | | Core hours | | |
|---|---|---|---|---|---|---|---|---|
| collocated benchmark type | Disaggregation | Ideal Non-sharing | Realistic | Disaggregation | Realistic | Disaggregation | Ideal Non-sharing | Realistic |
| BT, A | 0.937 | 0.893 | 0.693 | 0.877 | 1.0 | 0.968 | 1.0 | 1.29 |
| BT, W | 0.903 | 0.89 | 0.64 | 0.981 | 1.0 | 0.994 | 1.0 | 1.39 |
| CG, B | 0.992 | 0.901 | 0.65 | 0.94 | 1.0 | 0.908 | 1.0 | 1.39 |
| EP, B | 0.915 | 0.891 | 0.661 | 0.901 | 1.0 | 0.98 | 1.0 | 1.35 |
| LU, A | 0.941 | 0.893 | 0.674 | 0.929 | 1.0 | 0.964 | 1.0 | 1.33 |
| MG, A | 0.903 | 0.89 | 0.625 | 1.01 | 1.0 | 1.01 | 1.0 | 1.42 |
| MG, W | 0.903 | 0.89 | 0.638 | 1.01 | 1.0 | 1.0 | 1.0 | 1.39 |

**LULESH**

64 ranks, 2 nodes
32 out of 36 cores allocated.

**NAS**

1 – 4 ranks
Distributed across nodes.

# #1 CPU Sharing



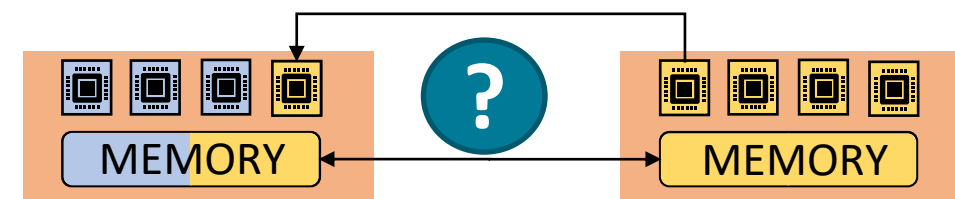| collocated benchmark type | Mean utilisation | | | Total time | | Core hours | | |
|---|---|---|---|---|---|---|---|---|
| | Disaggregation | Ideal Non-sharing | Realistic | Disaggregation | Realistic | Disaggregation | Ideal Non-sharing | Realistic |
| BT, A | 0.937 | 0.893 | 0.693 | 0.877 | 1.0 | 0.968 | 1.0 | 1.29 |
| BT, W | 0.903 | 0.89 | 0.64 | 0.981 | 1.0 | 0.994 | 1.0 | 1.39 |
| CG, B | 0.992 | 0.901 | 0.65 | 0.94 | 1.0 | 0.908 | 1.0 | 1.39 |
| EP, B | 0.915 | 0.891 | 0.661 | 0.901 | 1.0 | 0.98 | 1.0 | 1.35 |
| LU, A | 0.941 | 0.893 | 0.674 | 0.929 | 1.0 | 0.964 | 1.0 | 1.33 |
| MG, A | 0.903 | 0.89 | 0.625 | 1.01 | 1.0 | 1.01 | 1.0 | 1.42 |
| MG, W | 0.903 | 0.89 | 0.638 | 1.01 | 1.0 | 1.0 | 1.0 | 1.39 |

**LULESH**

64 ranks, 2 nodes
32 out of 36 cores allocated.

**NAS**

$1 - 4$ ranks
Distributed across nodes.

# #2 Serving Remote Memory



Slowdown of the batch job LULESH.

Overhead [%].

LULESH problem size
- 15
- 18
- 20
- 25

1 ms  5 ms  10 ms  25 ms  50 ms  100 ms  250 ms  500 ms  1 ms  5 ms  10 ms  25 ms  50 ms  100 ms  250 ms  500 ms

Slowdown of the batch job MILC.

Overhead [%].

MILC problem size
- 32
- 64
- 96
- 128

1 ms  5 ms  10 ms  25 ms  50 ms  100 ms  250 ms  500 ms  1 ms  5 ms  10 ms  25 ms  50 ms  100 ms  250 ms  500 ms

**LULESH, MILC –** 32 ranks, 1 node, 32 out of 36 cores allocated.

# #2 Serving Remote Memory



## Slowdown of the batch job LULESH.



LULESH problem size
- 15
- 18
- 20
- 25

## Slowdown of the batch job MILC.



MILC problem size
- 32
- 64
- 96
- 128

**LULESH, MILC –** 32 ranks, 1 node, 32 out of 36 cores allocated.

# #2 Serving Remote Memory



Slowdown of the batch job LULESH.

LULESH problem size: 15, 18, 20, 25

Slowdown of the batch job MILC.

Overhead below 7.5% when serving 1 GB/s.

MILC problem size: 32, 64, 96, 128

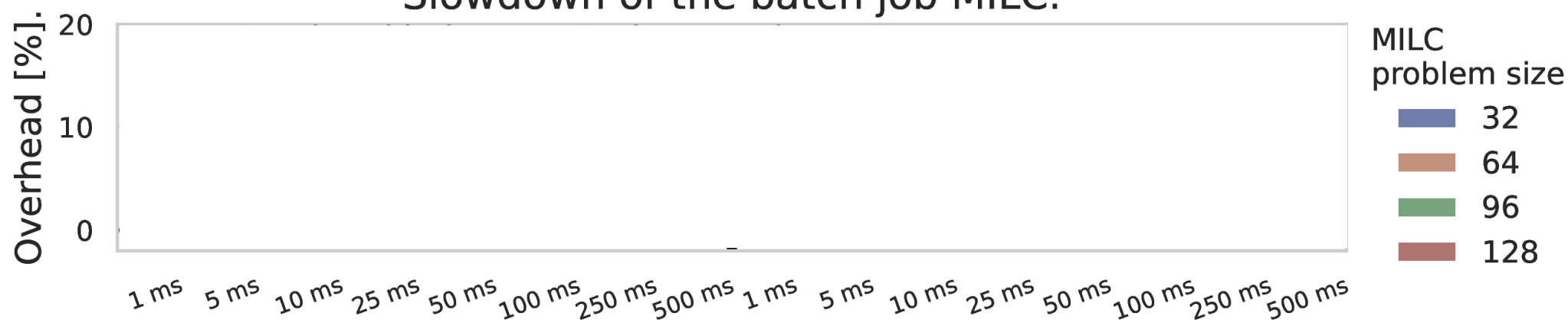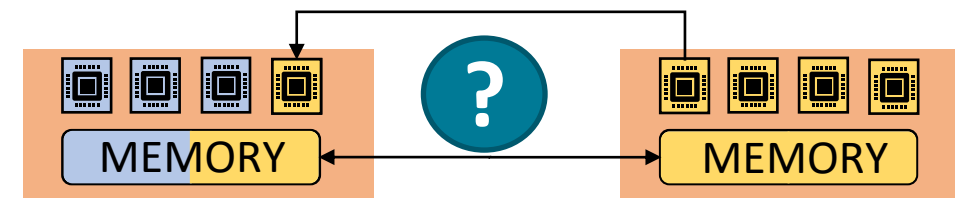**LULESH, MILC –** 32 ranks, 1 node, 32 out of 36 cores allocated.

# #2 Serving Remote Memory



Slowdown of the batch job LULESH.
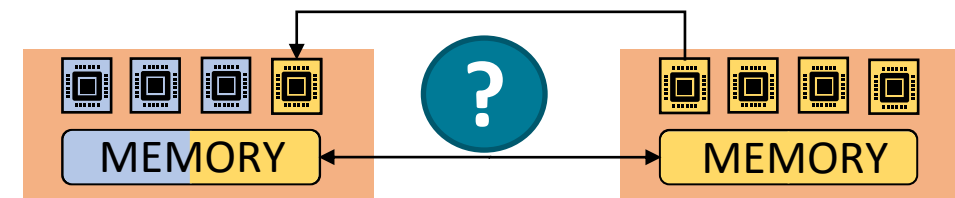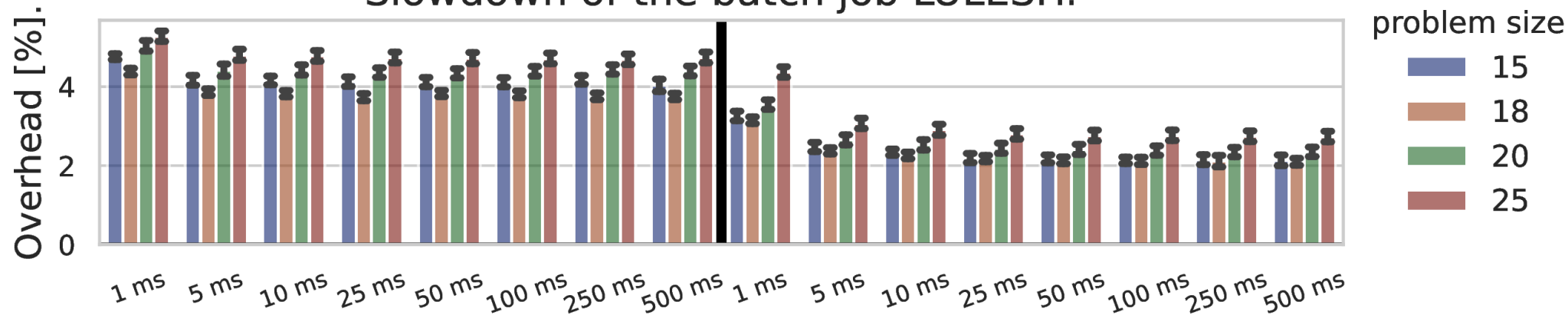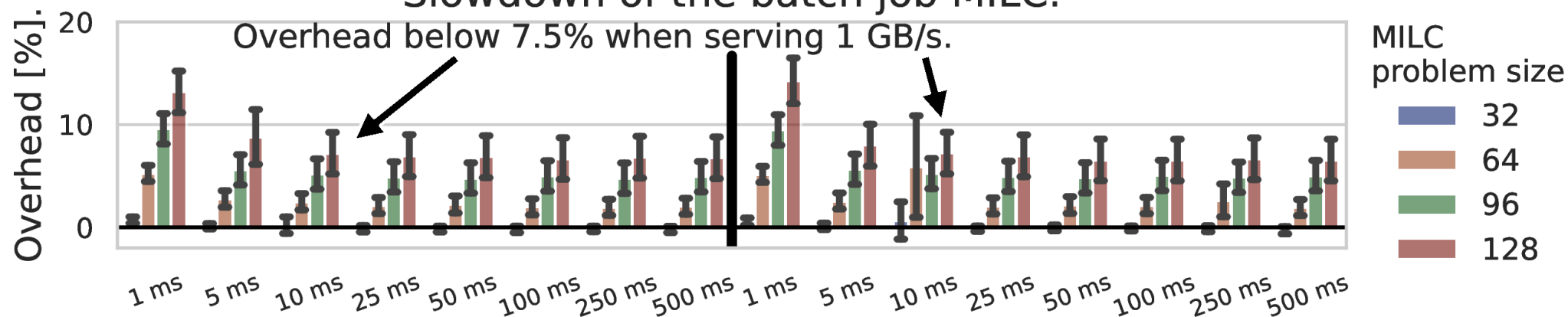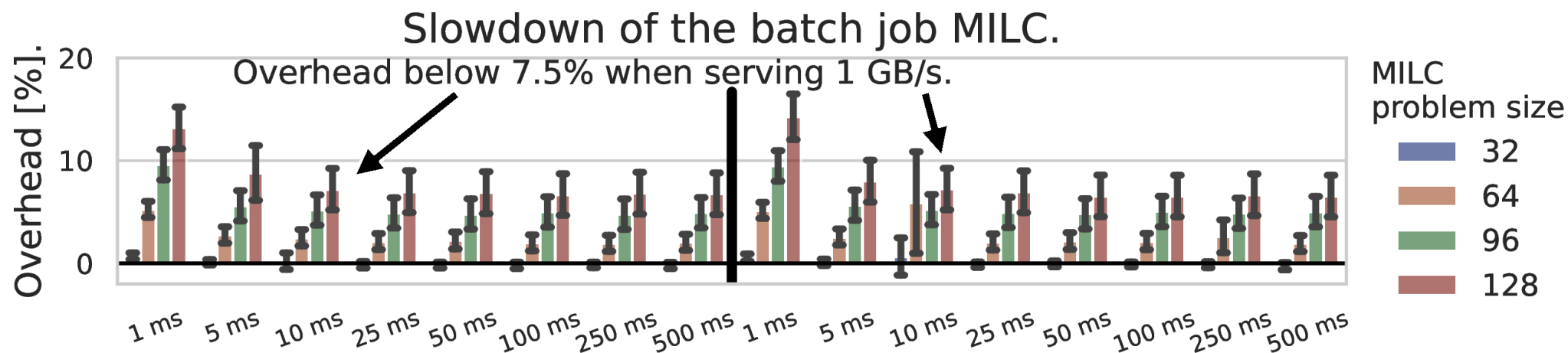
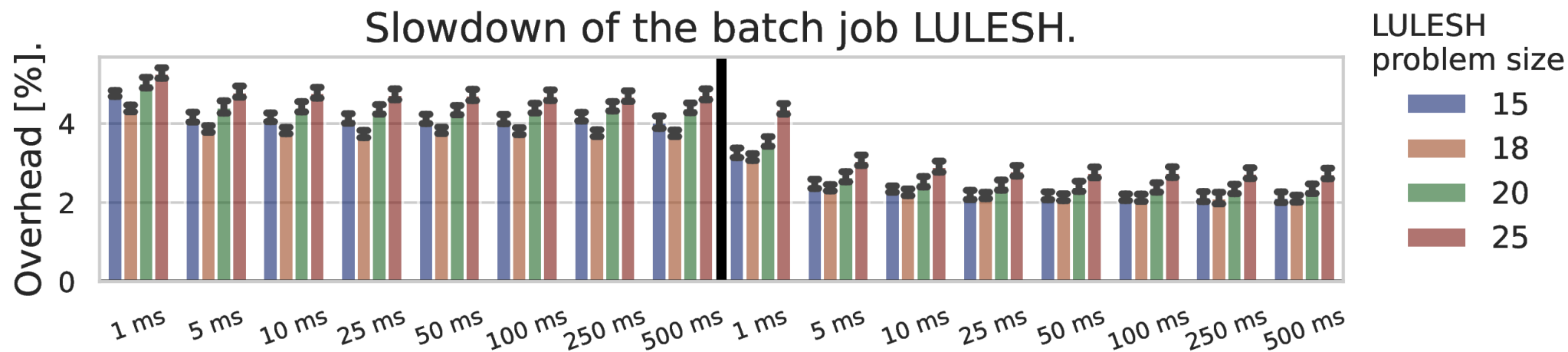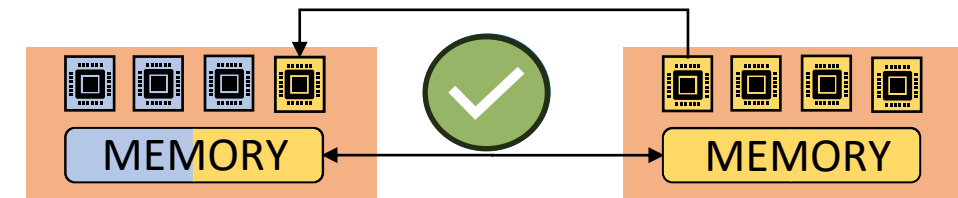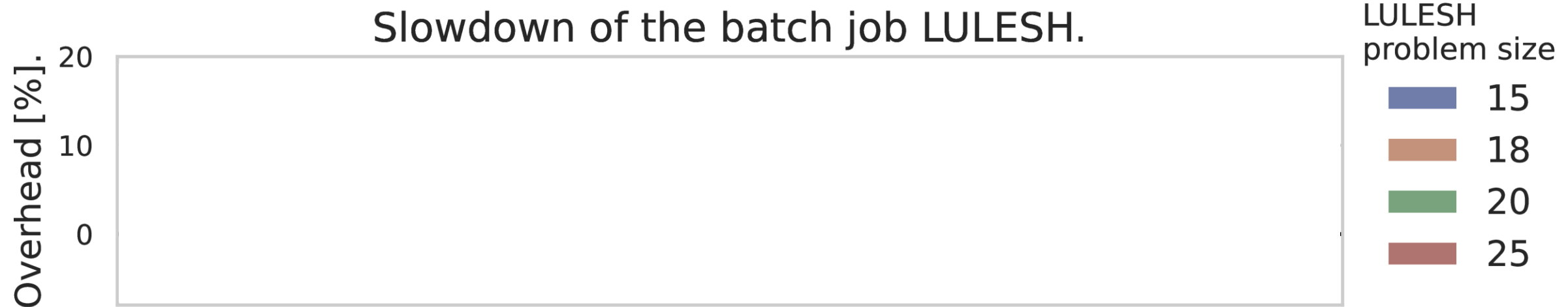Slowdown of the batch job MILC.

Overhead below 7.5% when serving 1 GB/s.

**LULESH, MILC –** 32 ranks, 1 node, 32 out of 36 cores allocated.

# #3 Co-locating GPU and CPU workloads



## Slowdown of the batch job LULESH.

Overhead [%].

20
10
0

Co-located GPU application.

LULESH problem size
- 15
- 18
- 20
- 25

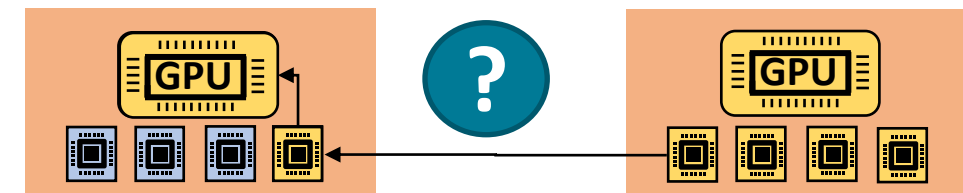**LULESH –** 27 ranks, 3 nodes, 9 out of 12 cores allocated.
**Rodinia –** 1 MPI rank, 1 GPU.

# #3 Co-locating GPU and CPU workloads

Slowdown of the batch job LULESH.

LULESH problem size
- 15
- 18
- 20
- 25

Overhead [%].

Baseline LULESH execution times: 24.5, 48.3, 74, 183.5 seconds.

Co-located GPU application.

bfs   gaussian   hotspot   myocyte   pathfinder   srad-v1

**LULESH –** 27 ranks, 3 nodes, 9 out of 12 cores allocated.
**Rodinia –** 1 MPI rank, 1 GPU.

# #3 Co-locating GPU and CPU workloads



Slowdown of the batch job LULESH.

Baseline LULESH execution times: 24.5, 48.3, 74, 183.5 seconds.

Co-located GPU application.

LULESH problem size: 15, 18, 20, 25

**LULESH –** 27 ranks, 3 nodes, 9 out of 12 cores allocated.
**Rodinia –** 1 MPI rank, 1 GPU.

# Summary



**HPC System Utilization - CPU**

80% and 70% of idle node events last less than 10 minutes.

# Summary



### HPC System Utilization - CPU

**Minimum**

**Maximum**

80% and 70% of idle node events last less than 10 minutes.



### Software Solution

**Standard HPC Node**

**Hardware Disaggregation**

✅ High performance
❌ Inflexible architecture

✅ High efficiency
❌ Cost, performance penalty

**Existing Coupled Hardware Systems**

**Software Abstraction for Disaggregation**

# Summary

# Summary



spcl/rFaaS

# Summary



HPC System Utilization - CPU

80% and 70% of idle node events last less than 10 minutes.



Software Solution



#1 CPU Sharing

## spcl/rFaaS

> *"the goal of achieving near 100% utilization while supporting a real parallel supercomputing workload is unrealistic"*

**Scheduling for Parallel Supercomputing:
A Historical Perspective of Achievable Utilization**

James Patton Jones[1] and Bill Nitzberg[1]

MRJ Technology Solutions
NASA Ames Research Center, M/S 258-6
Moffett Field, CA 94035-1000

jjones@nas.nasa.gov

1999