# Locality Aware Scheduling For Scalable Heterogeneous Environments

ROSS 2020

**Alok Kamatar, Ryan Friese, Roberto Gioiosa**

# Outline

- Motivation
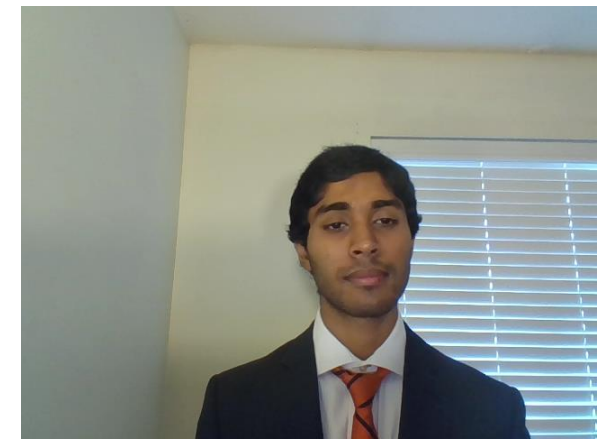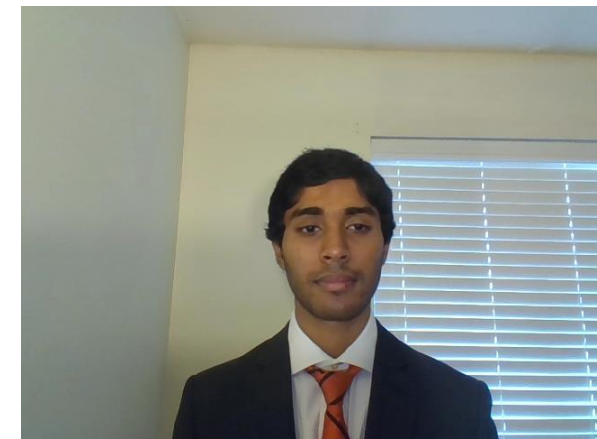
- The Minos Computing Library

- Architecture Details
  - Resident Memory
  - Scheduling Algorithm

- Experimental Results
  - SHOC Benchmarks
  - Hyper-parameter evaluation

- Conclusions and Future Work

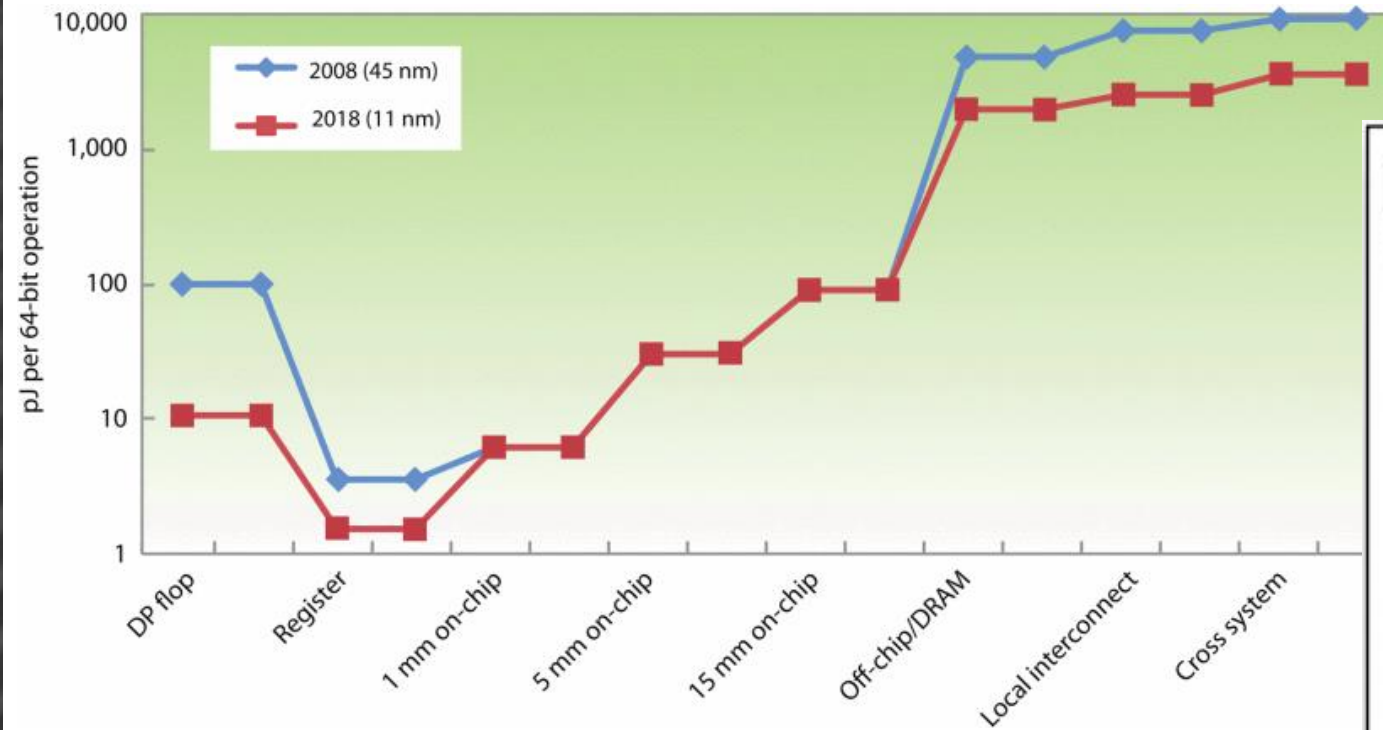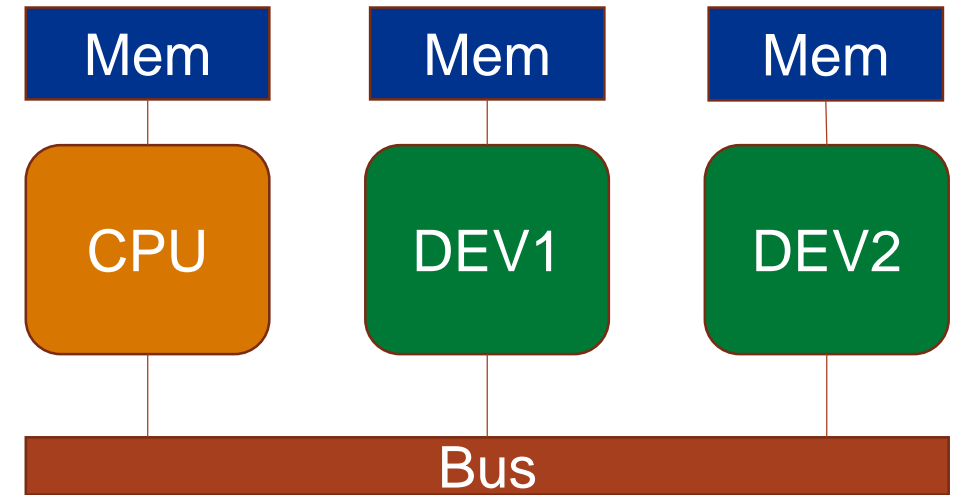# Exploiting Scalable, Heterogeneous Systems

- Increasing levels of hardware specialization
  - i.e. GPUs, Deep Learning Accelerators, DSPs, etc.
- Higher complexity to program new applications
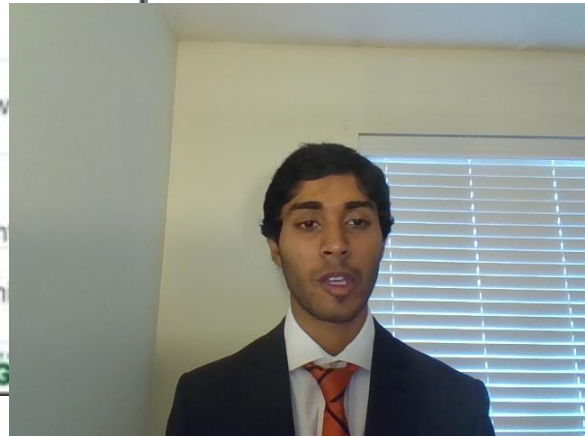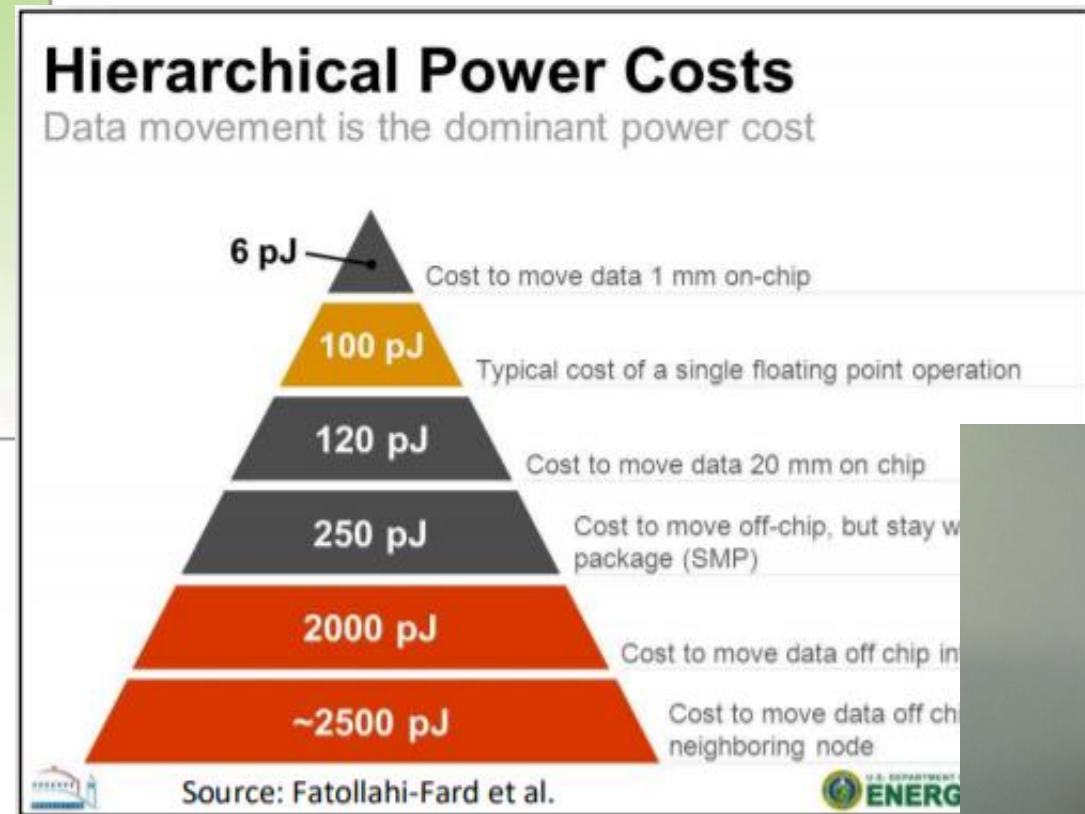- Potential pitfalls when scaling to other architectures

# Motivation

**Multiple Memory Nodes**

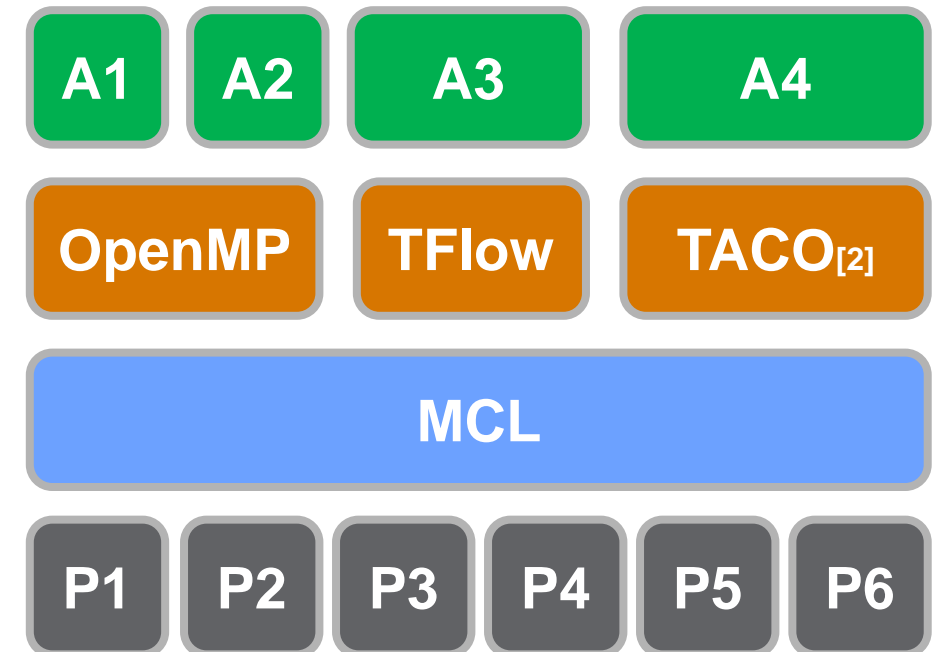- Data movement is now the critical aspect of performance
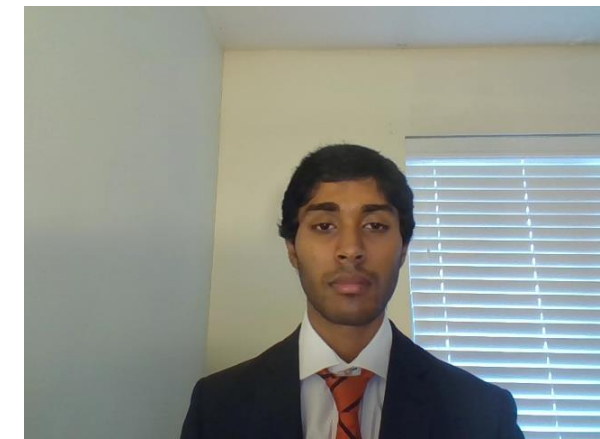


Source: Kogge, Shalf 2013

# The Minos Computing Library (MCL) [1]

- Framework for programming heterogeneous systems

- Features
  - Dynamic task scheduling onto available resources
  - Co-schedule independent applications
  - Internal profiling and tracing capabilities

- Flexibility
  - Automatic scaling to available resources
  - Works/can be integrated with commonly used technologies
  - Independent scheduler framework

| A1 | A2 | A3 | A4 |

| OpenMP | TFlow | TACO[2] |

| MCL |

| P1 | P2 | P3 | P4 | P5 | P6 |

[1] R. Gioiosa, B. Mutlu, S. Lee, J. Vetter, G. Picierro, M. Cesati. 2020. The Minos Computing Library: Efficient Parallel Programming for Extremely Heterogeneous Systems. In General Purpose Processing Using GPU (GPGPU '20), February 23, 2020, San Diego, CA, USA
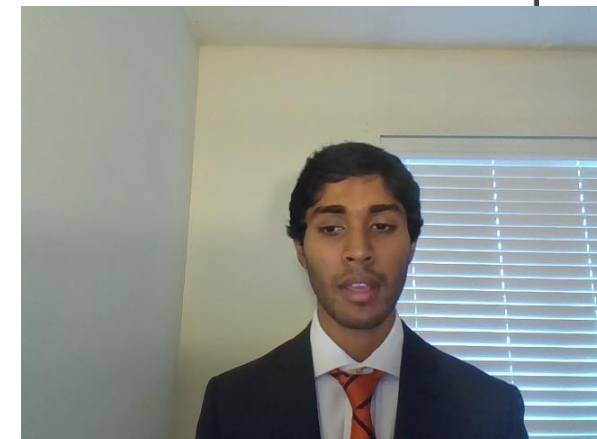[2] G. Kestor, R. Gioiosa, M. Raugas. Towards Performance Portability through an Integrated Programming Eco-System for Tensor Algebra. In Performance, Portability, and Productivity in HPC Forum, to be held online, September 2020

# MCL Program Example

- Same piece of code will exploit all available resources

- OpenCL kernels are directly usable by MCL

```
1    #include <minos.h>
2
3    int main(void){
4        mcl_init(NWorkers, <flags>);
5    ....
6
7        for(i=0; i< NIter; i++){
8            hdl[i] = mcl_task_create(<filename>, <kernel>, <flags>);
9            mcl_task_set_arg(hdl[i], 0, <addr>, <size>, <flags>);
10           ...
11           mcl_task_set_art(hdl[i], k, <addr>, <size>, <flags>);
12           mcl_exec(hdl[i], <pes>, NULL, MCL_TASK_GPU);
13       }
14   ...
15       mcl_wait_all();
16   ...
17       mcl_finit();
18
19       return 0;
20   }
21
```
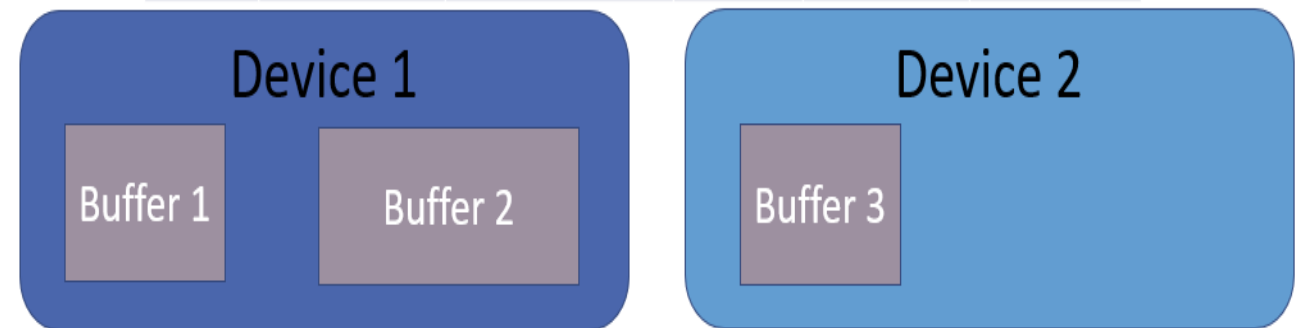
# MCL Resident Memory Module

- Allows persistent data to remain in device memory across tasks

- Coordinates data movement so correct data is transferred to the correct device

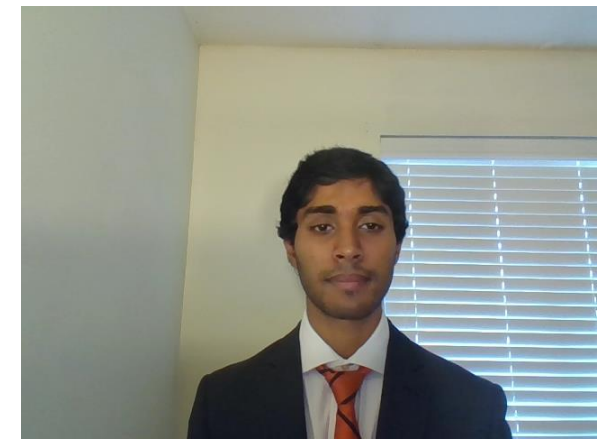- Supports read-only (i.e. multiple copies) and read-write data (exclusive copies)

| Pid | Memory | Concurrent Uses | Size | Device 1 | Device 2 |
|-----|--------|-----------------|-------|----------|----------|
| 1 | 1 | 1 | 10 MB | 1 | 0 |
| 1 | 2 | 1 | 20 MB | 1 | 0 |
| 1 | 3 | 1 | 10MB | 0 | 1 |
| ... | ... | ... | ... | ... | ... |

**Device 1**
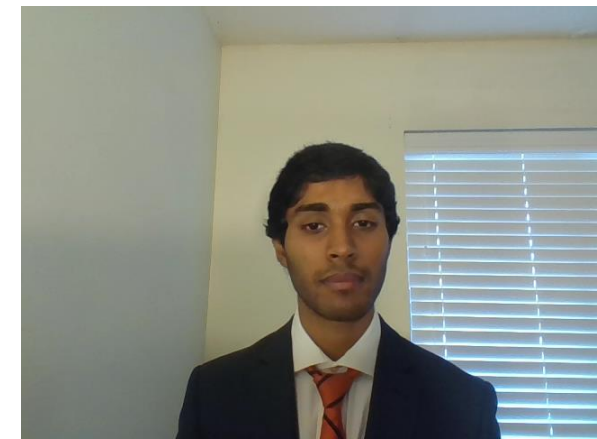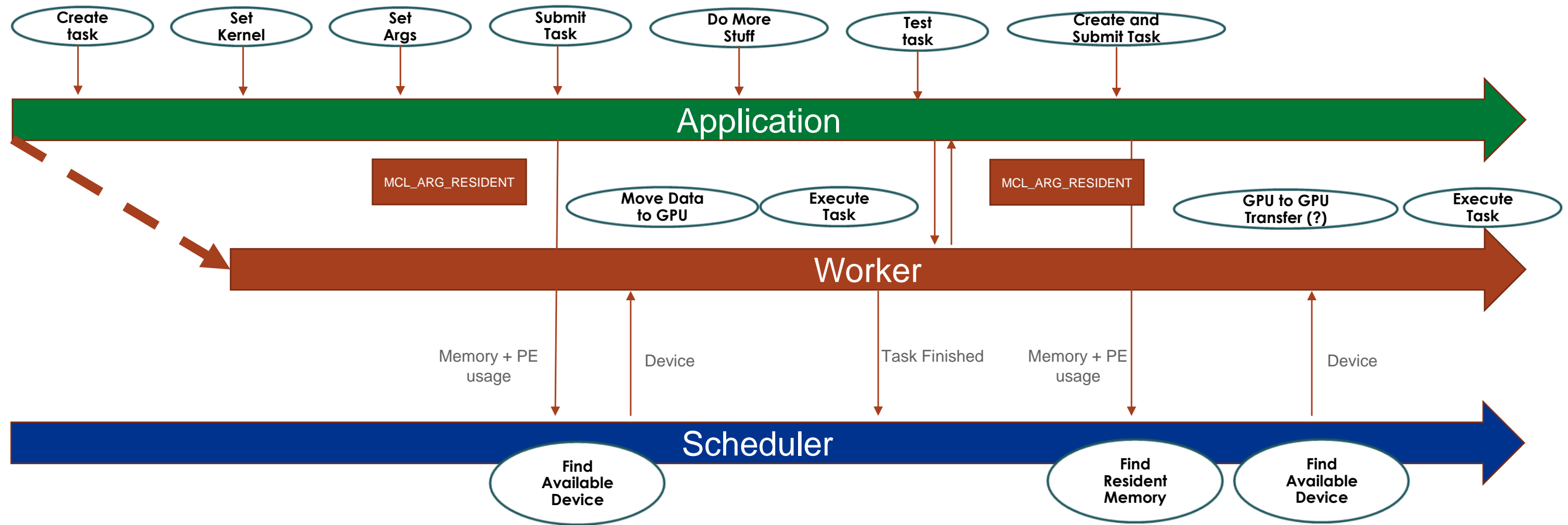
Buffer 1    Buffer 2

**Device 2**

Buffer 3

MCL_ARG_RESIDENT
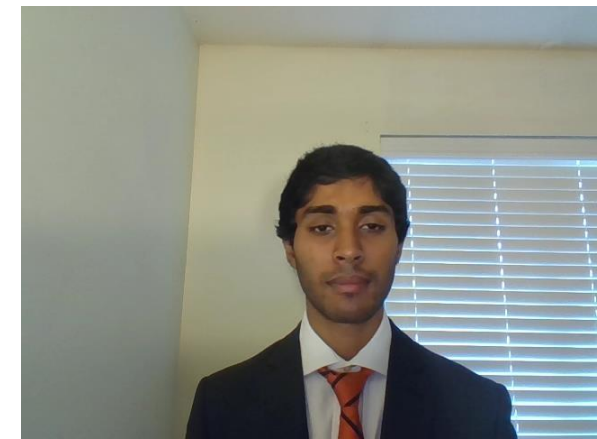
MCL_ARG_INVALID

MCL_ARG_DONE

# MCL Trace with Resident Memory

# MCL Schedulers: Round-Robin Scheduler

- Schedules tasks in a first-in first out manner to the next available device

- Typically achieves high resource utilization

- Problems
  - Suffers from head-of-line blocking
  - Makes poor scheduling decisions when there is data reuse between kernels
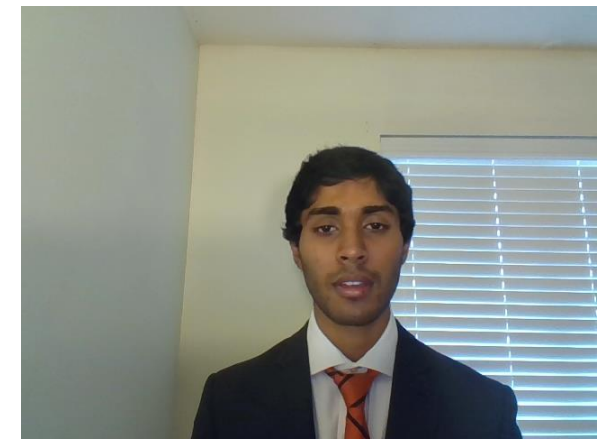  - **Unnecessary data transfers**
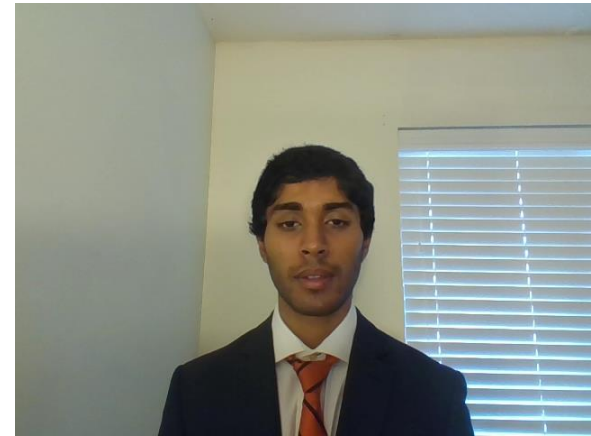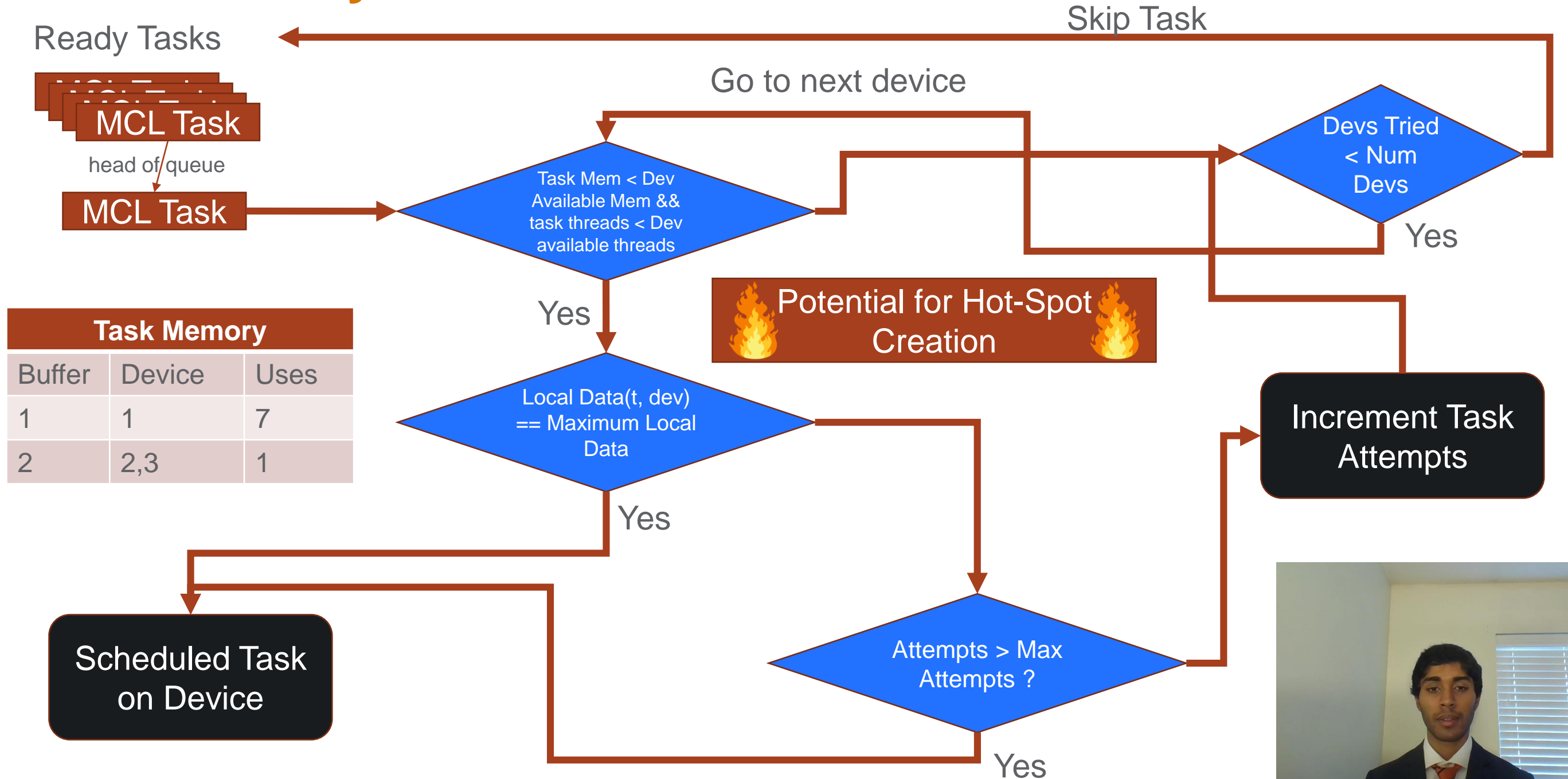
# MCL Schedulers: Delay Scheduler

- Delays kernels from running on devices without device local data to minimize data transfers

- Skips devices that do not have device local data

- Skips tasks when waiting for busy devices

- Limits the number of times a task can be delayed to prevent a task from blocking too long

---

**Algorithm 1** Device Local Data

1: **LocalData(Device $\delta$, Task $t$):**
2: $bytes \leftarrow 0$
3: **for** all $\beta$ in $t.buffers$ **do**
4:    **if** $\beta$ in $\delta.data$ **then**
5:       $bytes \leftarrow bytes + \beta.size$
6:    **end if**
7: **end for**
8: return($bytes$)
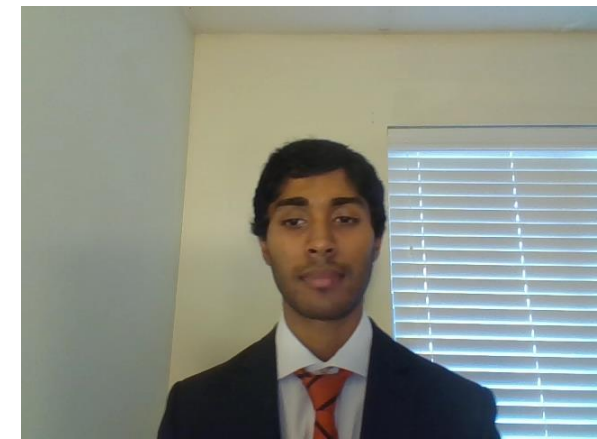9: **end LocalData**

# Delay Scheduler Illustration

# MCL Schedulers: Mixed Scheduler

- Attempts to balance system utilization with data locality concerns

- Detects popular pieces of data to create replicas

- The best predictor of how popular data will be is how popular it was in the past

---
**Algorithm 3** Device Local Data With Copy Factor
---

1: **LocalDataCopyFactor(Device** $\delta$, **Task** $t$, **CopyFactor** $\gamma$)**:**
2: $bytes \leftarrow 0$
3: **for** all $\beta$ in $t.buffers$ **do**
4:    **if** $\beta$ in $\delta.data$ **and** DEVICES($\beta$) $\geq$ $\log_\gamma$ TASKS($\beta$) **then**
5:       $bytes \leftarrow bytes + \beta.size$
6:    **end if**
7: **end for**
8: return($bytes$)
9: **end LocalDataCopyFactor**
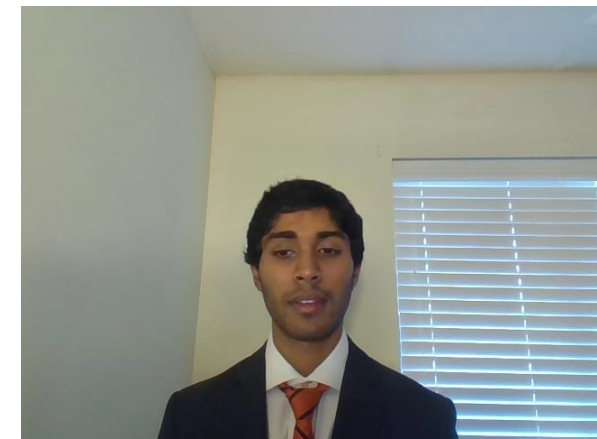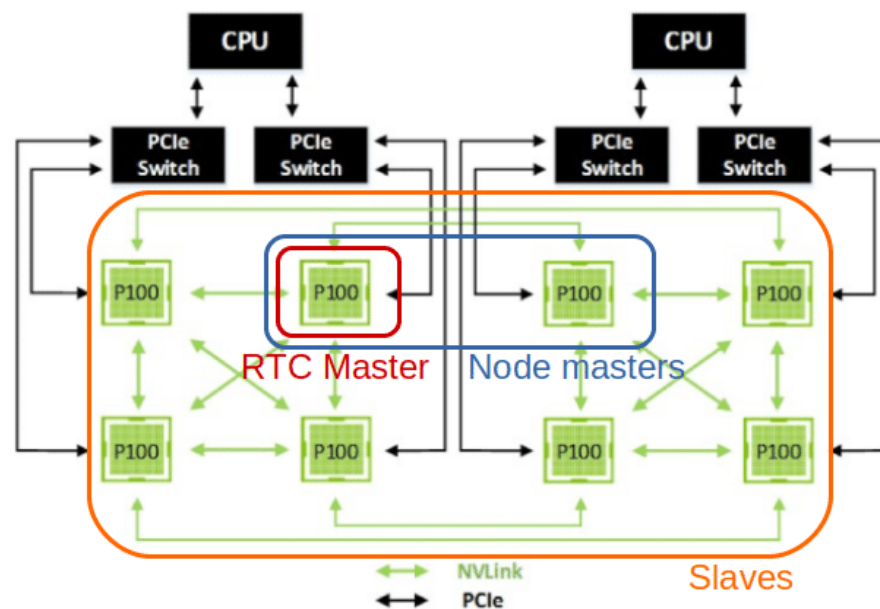
# Experimental Setup

- Benchmarks (from SHOC)
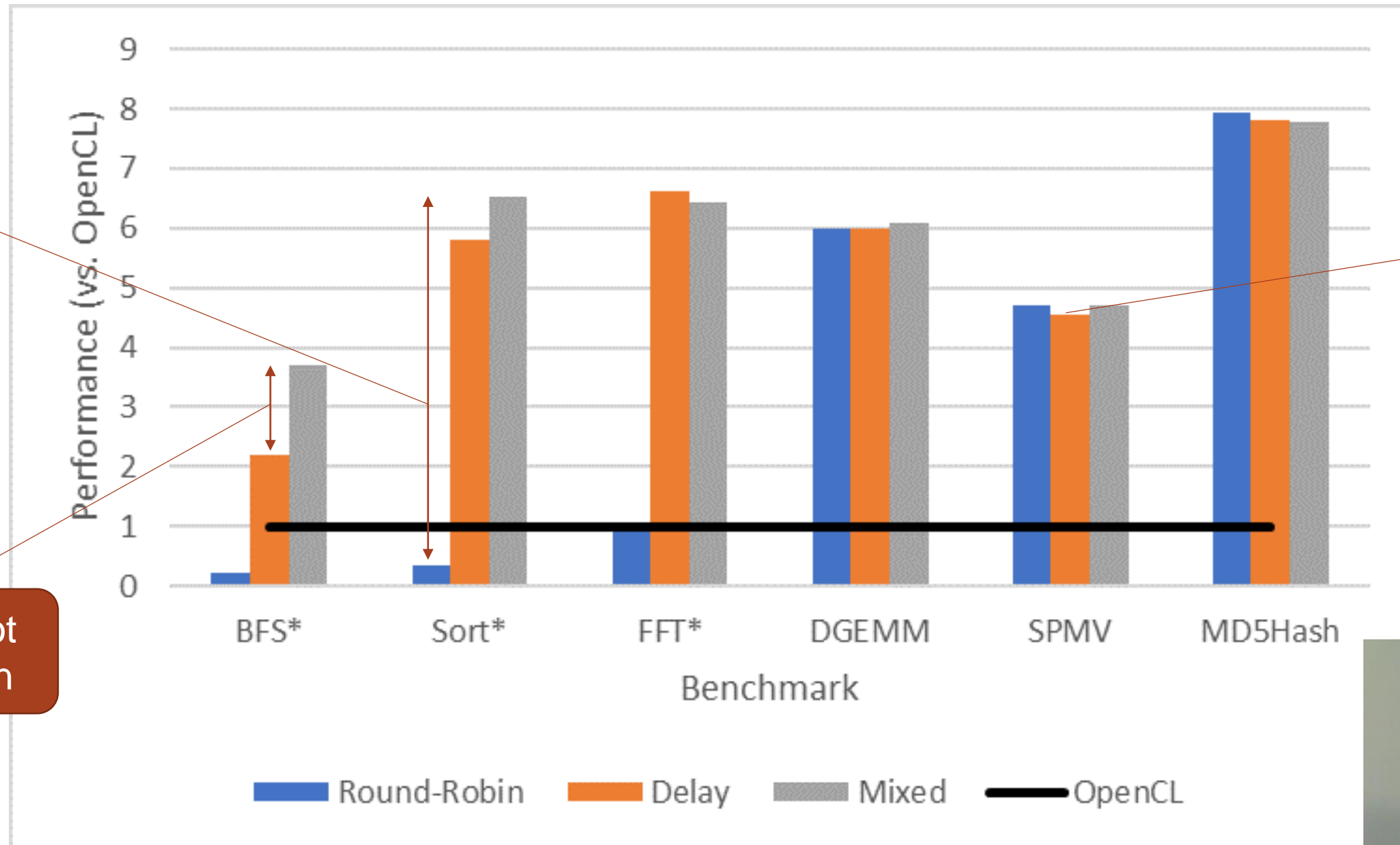  - BFS*
  - Sort*
  - FFT*
  - DGEMM
  - SPMV
  - MD5Hash

* = Exploits Data Locality

**DGX-1 P100**
- 2x 20 Core Intel Xeon CPU
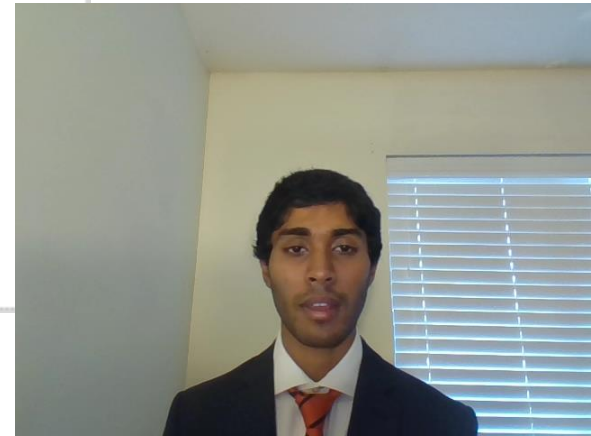- 8x Nvidia P100 GPU, NVLINK
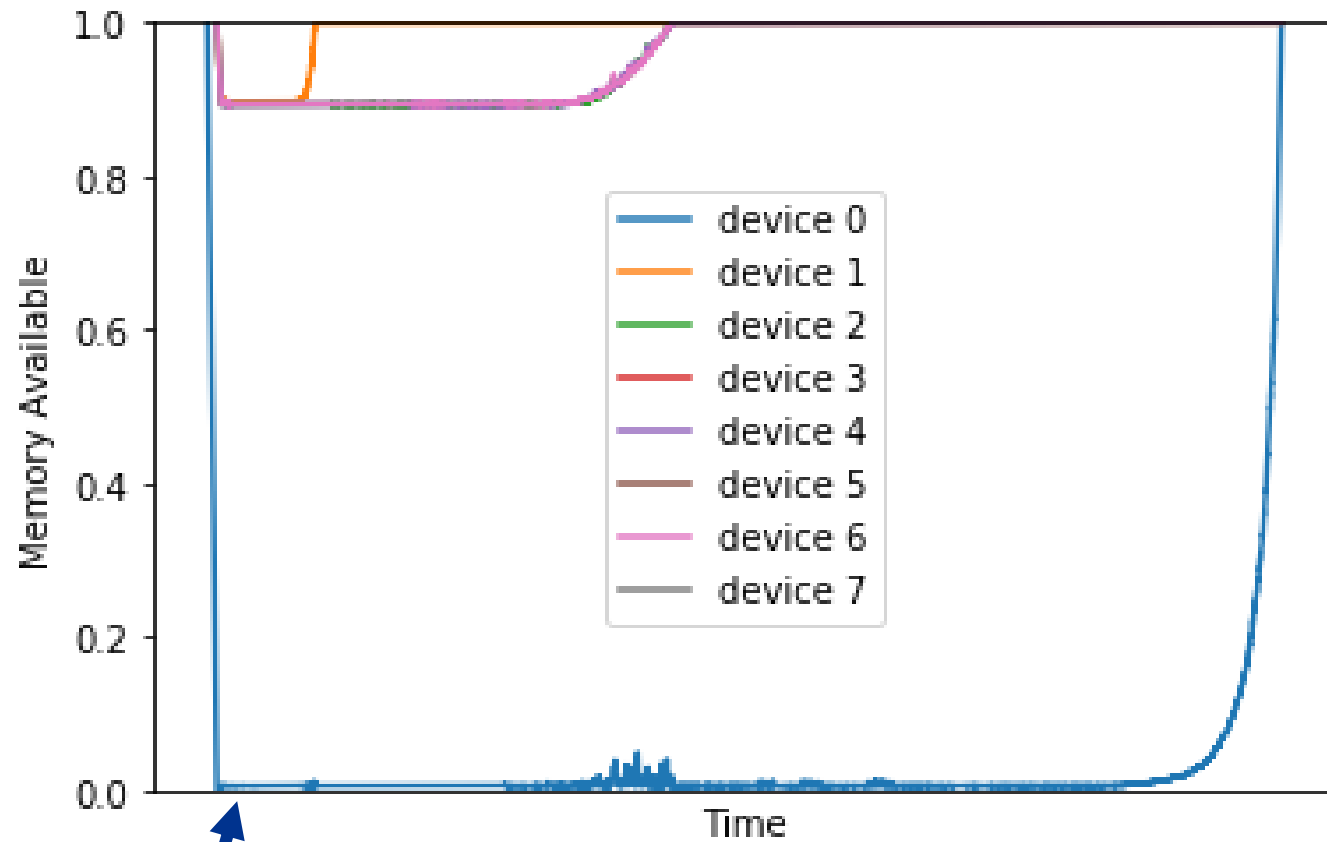- 256 GB RAM + 16 GPU

# Scheduling Evaluation

# BFS Memory Usage Comparison

Delay Scheduler Memory Trace

Mixed Scheduler Memory Trace



Hot Spot creation on device 1

GPU 1

Adaptive replication of popular data

Memory use from MCL internal trace capabilities

# Effect of Hyperparameters on Performance



BFS MTEPS (y-axis)
Max Attempts (x-axis): 2, 4, 8, 16, 32

Legend: Copy Factor 2, Copy Factor 4, Copy Factor 8, Copy Factor 16

**BFS Benchmark**
- 1,000,000 vertices
- 4096 Tasks

Num GPUs

When Max Attempts < Num Devices – low performance

If Max Attempts > Num Devices – copy factor dominates performance

# Effect of Hyperparameters on GBs Transferred
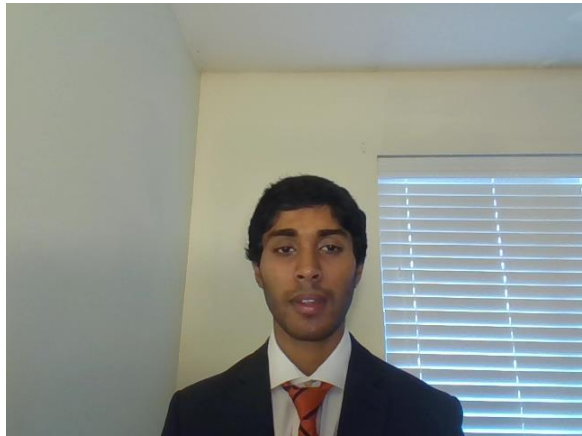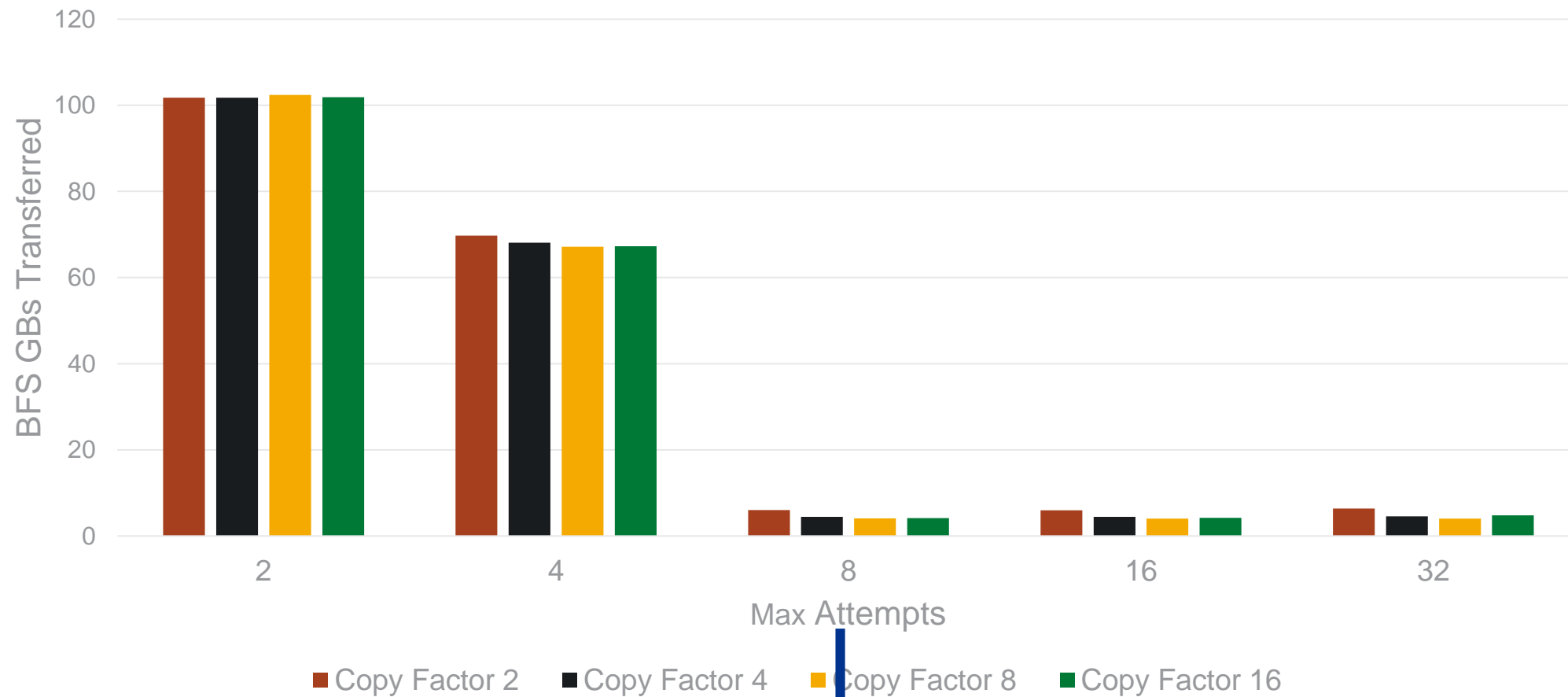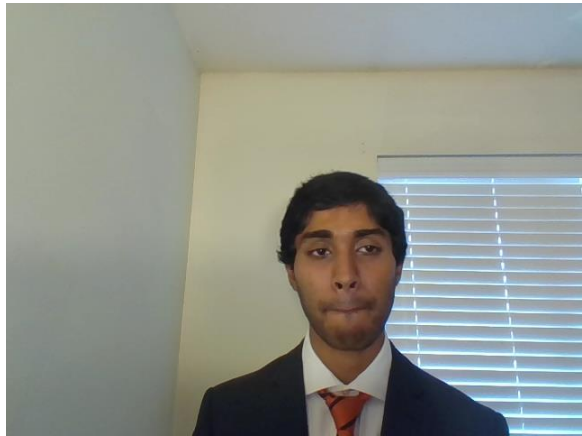


BFS Benchmark
- 1,000,000 vertices
- 4096 Tasks

When Max Attempts < Num Devices – low performance
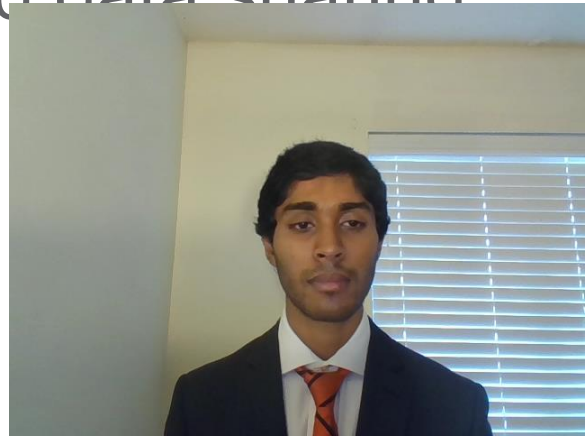
Num GPUs

If Max Attempts > Num Devices – copy factor dominates performance

# Conclusions + Future Work

- One of the primary challenges of using multiple devices is managing memory and coordinating data movements

- We introduce MCL Resident Memory to seamlessly manage device memory and coordinate data transfers

- Results demonstrate the importance of locality in achieving performance speedups

- Improves performance against other baselines

- Improve scheduler by dynamically determining hyperparameters

- In the future we want to expand this work to allow scheduling and data sharing among multiple applications

# Thank you

Alok Kamatar
Alok.Kamatar@pnnl.gov
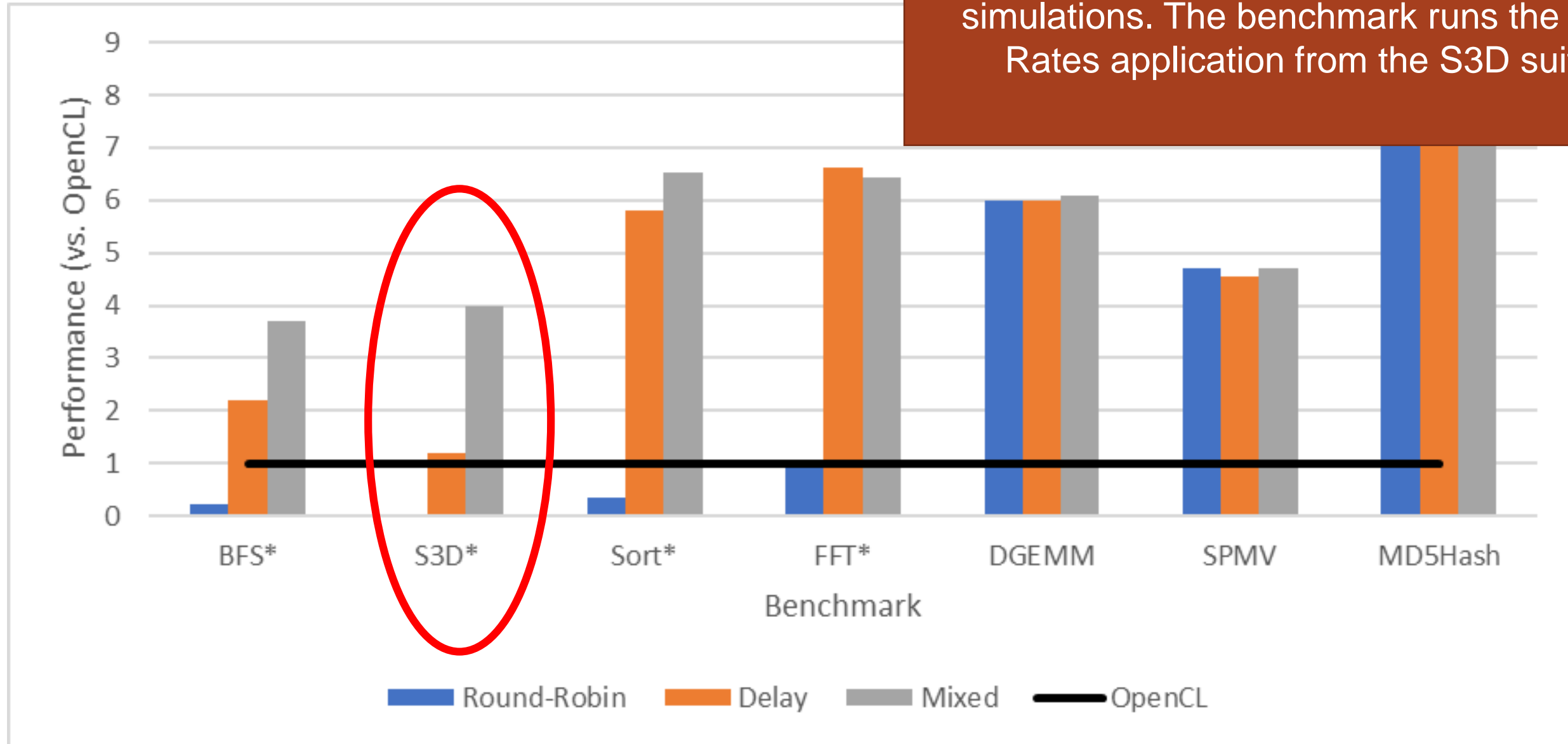
Roberto Gioiosa
Roberto.Gioiosa@pnnl.gov

Additional Slides

# Results With S3D

S3D is a standard direct numerical solver (DNS) used for combustion and other chemical simulations. The benchmark runs the Get-Rates application from the S3D suite
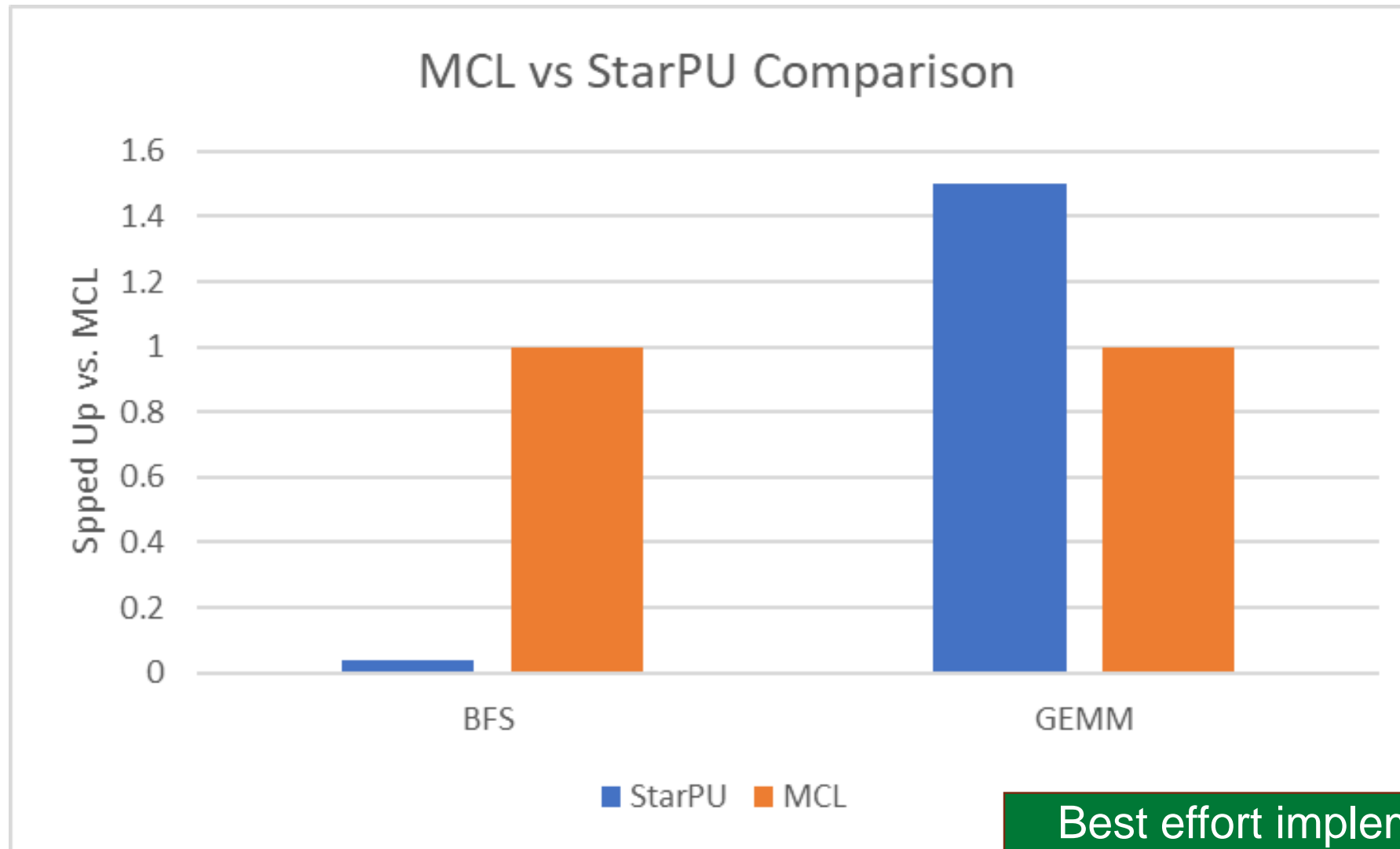
# Application Composition



**Experiment Details:**
- OpenCL application partitioned among 8 devices
- A variable mix of small (64x64) and large (1024x1024) GEMMS

# **Eviction/Checkpointing**

- Memory Usage is a limited resource that is under demand in a HPC system

- MCL supports flexible eviction policies that can be combined with scheduler policies

- When applications are unable to be run because no device has enough available memory, resident data can be evicted back to main memory

- To the user, MCL still behaves the same

- Currently supports a LRU policy

# Comparison Against StarPU



MCL vs StarPU Comparison

**StarPU Scheduling LWS:**
- Task is automatically scheduled on worker that released it
- Idle workers use data transfer performance estimates to determine weather to move data and "steal" work

Best effort implementation of StarPU – still investigating performance discrepancy