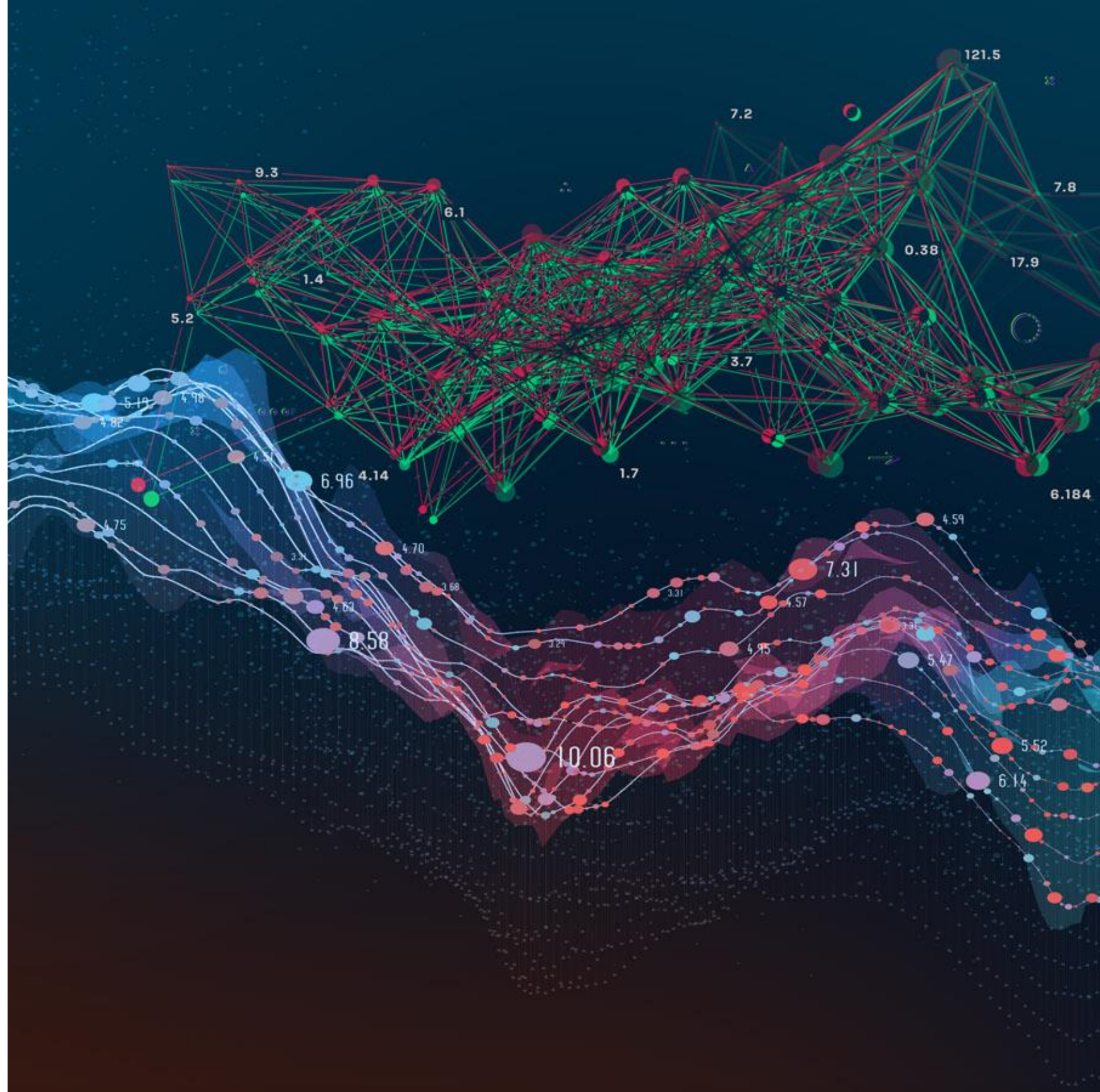


# MCL Multi-GPU Support + Schedulers

PPoPP '21

Alok Kamatar, Ryan Friese, Roberto Gioiosa



# Use Cases

- Translating OpenCL code from one device to multiple:

```
For dev in devs:  
    clCreateContext ...  
    clCreateProgramWithSrc ...  
    clCreateCommandQueueWithProps ..  
    clCreateKernel ...  
    //Create a buffer for each array  
    //etc.
```

- Lots of complexity overhead for each device
- Coordinate which devices are busy and idle
- Coordinate memory between devices



## Use Cases

- Translating MCL Code from single device to multiple devices

```
mcl_exec(task, pes, block_size, MCL_TASK_GPU);
```



```
mcl_exec(task, pes, block_size, MCL_TASK_GPU);
```

- Same code runs on all available resources
- MCL takes care of and hides complexity overhead

# Example – Easy Scaling

Select kama097@cenatehub:~

0	Tesla P100-SXM2...	On	00000000:06:00.0	Off	0	
N/A	45C	P0	49W / 300W	1327MiB / 16280MiB	100%	Default

1	Tesla P100-SXM2...	On	00000000:07:00.0	Off	0	
N/A	42C	P0	34W / 300W	10MiB / 16280MiB	0%	Default

2	Tesla P100-SXM2...	On	00000000:0A:00.0	Off	0	
N/A	40C	P0	32W / 300W	10MiB / 16280MiB	0%	Default

3	Tesla P100-SXM2...	On	00000000:0B:00.0	Off	0	
N/A	32C	P0	31W / 300W	10MiB / 16280MiB	0%	Default

4	Tesla P100-SXM2...	On	00000000:85:00.0	Off	0	
N/A	33C	P0	31W / 300W	10MiB / 16280MiB	0%	Default

5	Tesla P100-SXM2...	On	00000000:86:00.0	Off	0	
N/A	38C	P0	33W / 300W	10MiB / 16280MiB	0%	Default

6	Tesla P100-SXM2...	On	00000000:89:00.0	Off	0	
N/A	37C	P0	33W / 300W	10MiB / 16280MiB	0%	Default


7	Tesla P100-SXM2...	On	00000000:8A:00.0	Off	0	
N/A	39C	P0	32W / 300W	10MiB / 16280MiB	0%	Default


Processes:					GPU Memory
GPU	PID	Type	Process name		Usage
0	44449	C	./procBin_2_all		1317MiB

\$

kama097@cenatehub:~

\$

  
Spotify

RECORDED WITH  
SCREENCAST  MATIC

## Running Commands

- To replicate the results in the video, you can run these commands on your own GPU enabled system

```
>>> mcl_sched -s fffs -p rr &
>>> cd mcl/test
>>> make mcl_gemm
>>> ./mcl_gemm -r 8092 -w 16
-----
>>> export CUDA_VISIBLE_DEVICES=1,2,3,4,5,6,7
>>> mcl_sched -s fffs -p rr &
>>> cd mcl/test
>>> make mcl_gemm
>>> ./mcl_gemm -r 8092 -w 16
```

Commands used for DGX-1 demo 1

# MCL Schedulers

- MCL scheduler is a separate process that determines where and when tasks are run
- Flexible – i.e. different scheduler policies divided into two classes
  - First In First Out – The first task in the queue must be scheduled first,
  - FFFS scheduler – Allows the resource scheduler to select and of the currently enqueued tasks
- Different Policies
  - First Fit
  - Round Robbin
  - Delay Scheduler
  - Hybrid Scheduler

```
$ ~/local/bin/mcl_sched -h
Usage: /home/kama097/local/bin/mcl_sched [options]
       -s, --sched-class {fifo|fffs}  Select scheduler class (def = 'fifo')
       -p, --res-policy {ff|rr|delay|hybrid}  Select resource policy (def = class dependant)
       -h, --help                      Show this help
$ _
```



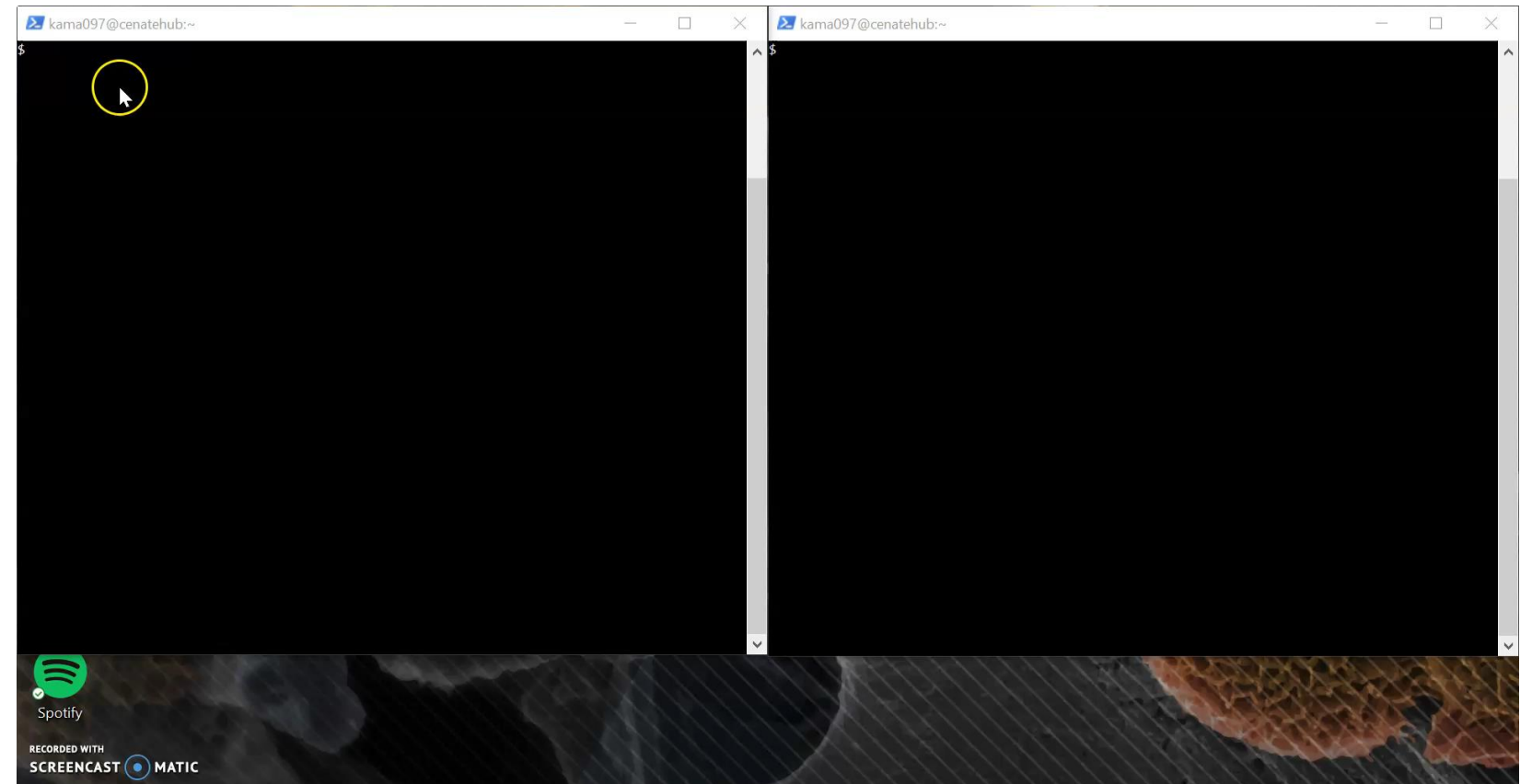
# First Fit Scheduler

- Schedules the next available task onto the current device if there is space
- Only moves onto the next device if the current device is full
- Space is determined by both memory and available processing threads
- Low Overhead

```
>>> mcl_sched -s fffs -p ff &
```

# To Run First Fit Scheduler with Trace

```
>>> make clean
>>> ./configure --enable-trace
>>> make && make install
>>> mcl_sched -s fffs -p ff &
>>> cd test
>>> make mcl_gemm
>>> ./mcl_gemm -r 8092 -w 16
```





\$



^

\$

^

\$



Spotify

RECORDED WITH  
**SCREENCAST**  **MATIC**

# First Fit Scheduler Results

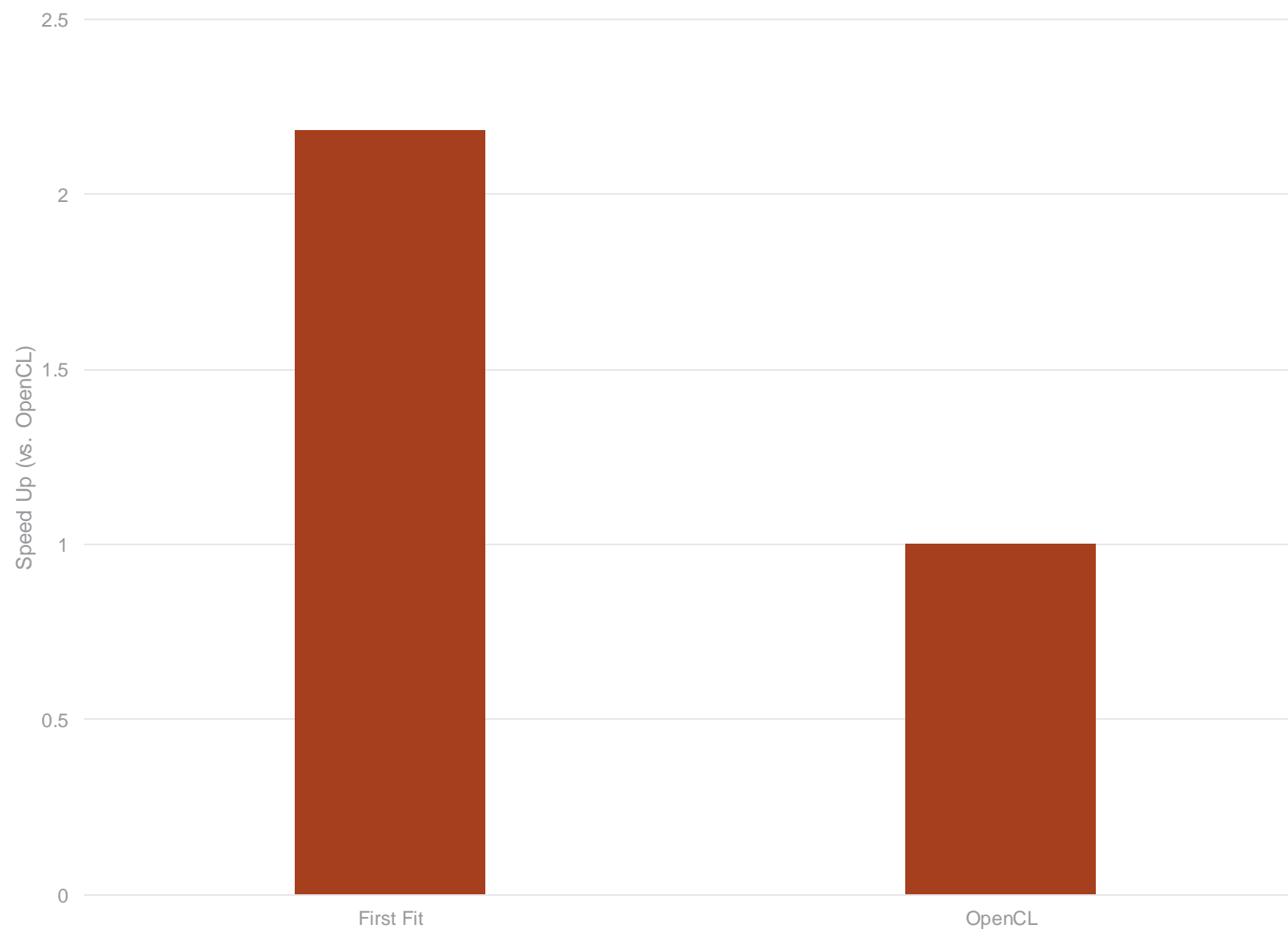
## DGEMM Benchmark:

- 64 X 64 matrices
- 1024 tasks

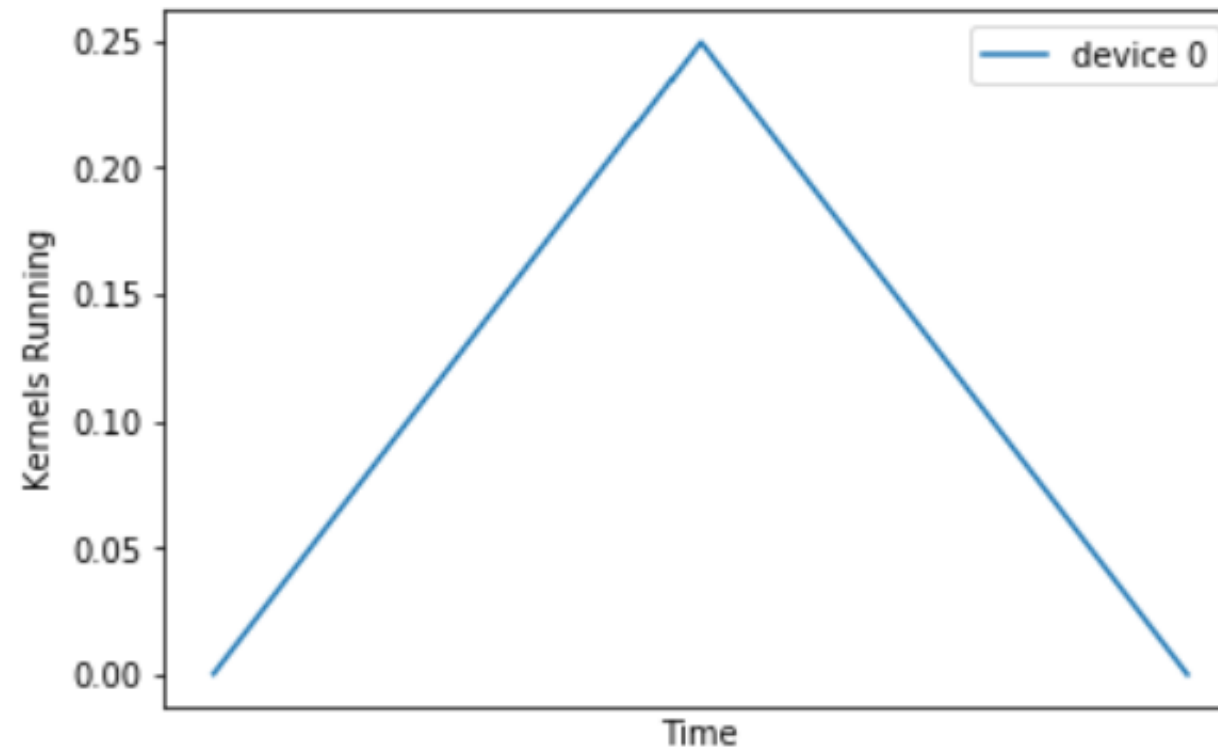
## Nvidia DGX-1 System:

- 2 Intel Xeon CPUs
- 8 Nvidia Pascal GPUS
- 256 GB memory + 16 GB per GPU memory

MCL Scheduler Comparison



# Scheduler Trace



- MCL scheduler trace can be enable by compiling with `-enable-trace`
- Allows us to visualize memory use, processing element use, and number of busy tasks on each device
- Where are the rest of the GPUS?
  - Some GPUs are heavily utilized; rest of the GPUs are under utilized



# Round Robin Scheduler

- Schedules tasks to devices in a circular manner
- Maintains a queue of devices.
  - Each incoming task gets assigned to the next device in the queue that is compatible
  - That device is moved to the back of the queue
- Typically achieves good full-system utilization

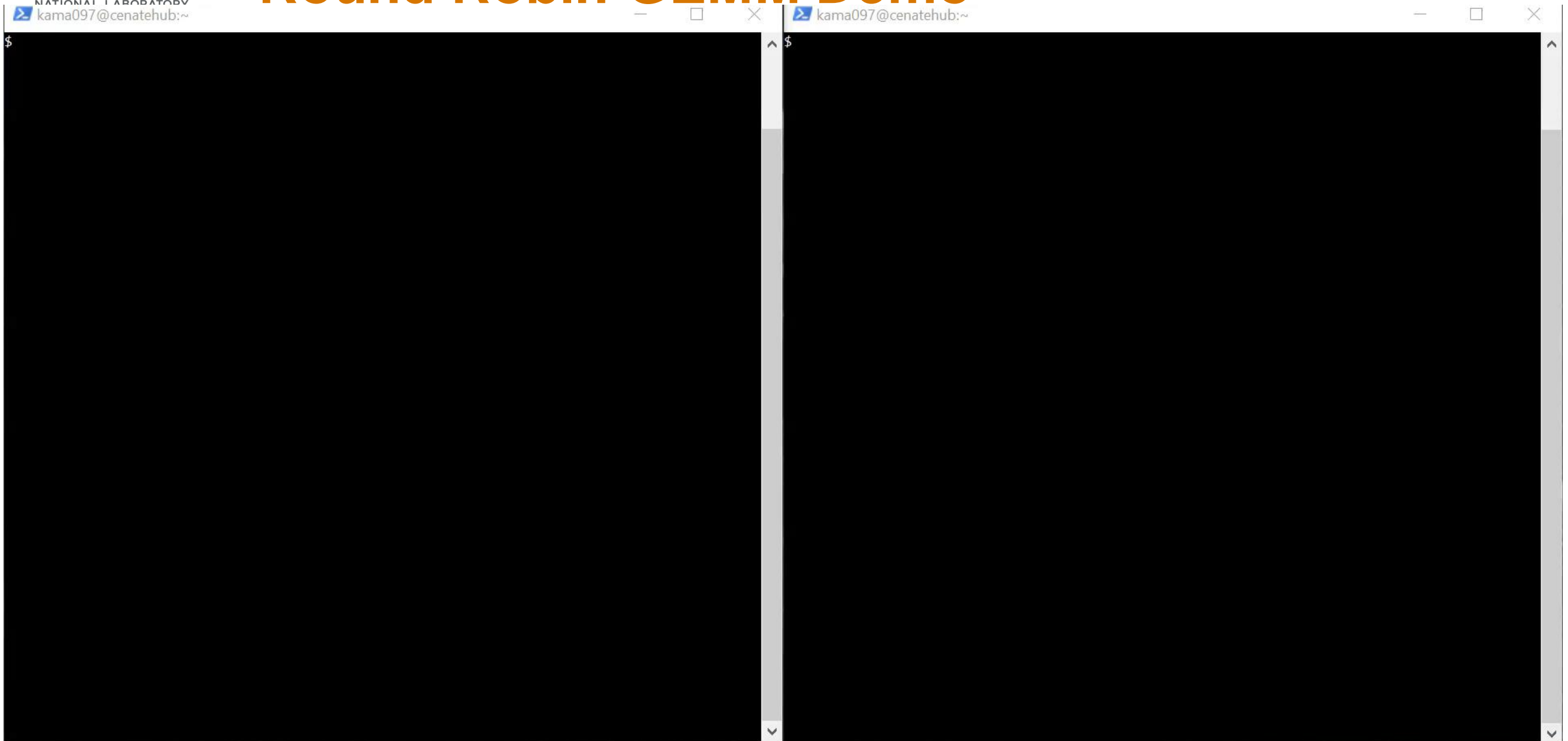
---

```
>>> mcl_sched -s fffs -p rr &
```



Pacific  
Northwest  
NATIONAL LABORATORY  
kama097@cenatehub:~

# Round Robin GEMM Demo



# Round Robin Scheduler Results

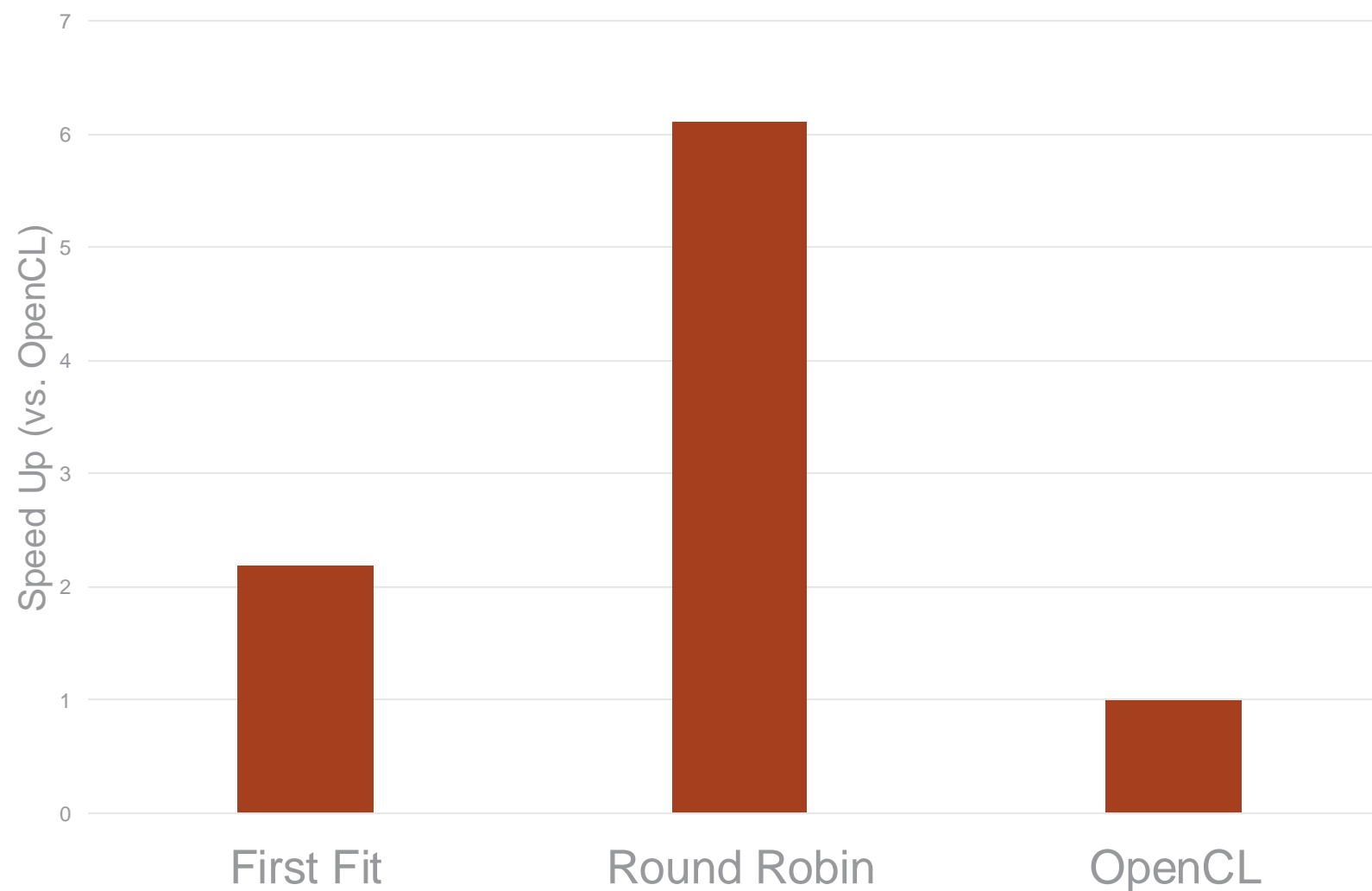
## DGEMM Benchmark:

- 64 X 64 matrices
- 1024 tasks

## Nvidia DGX-1 System:

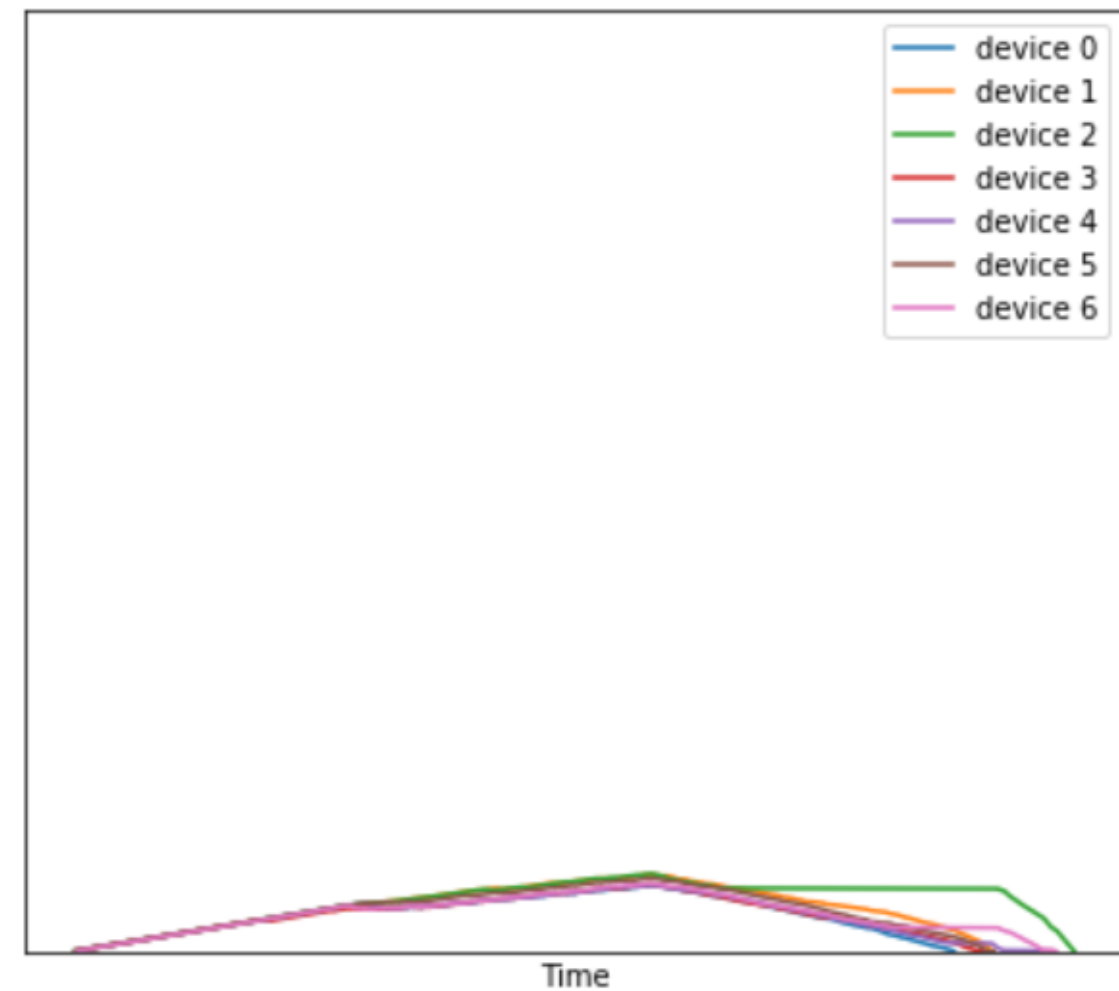
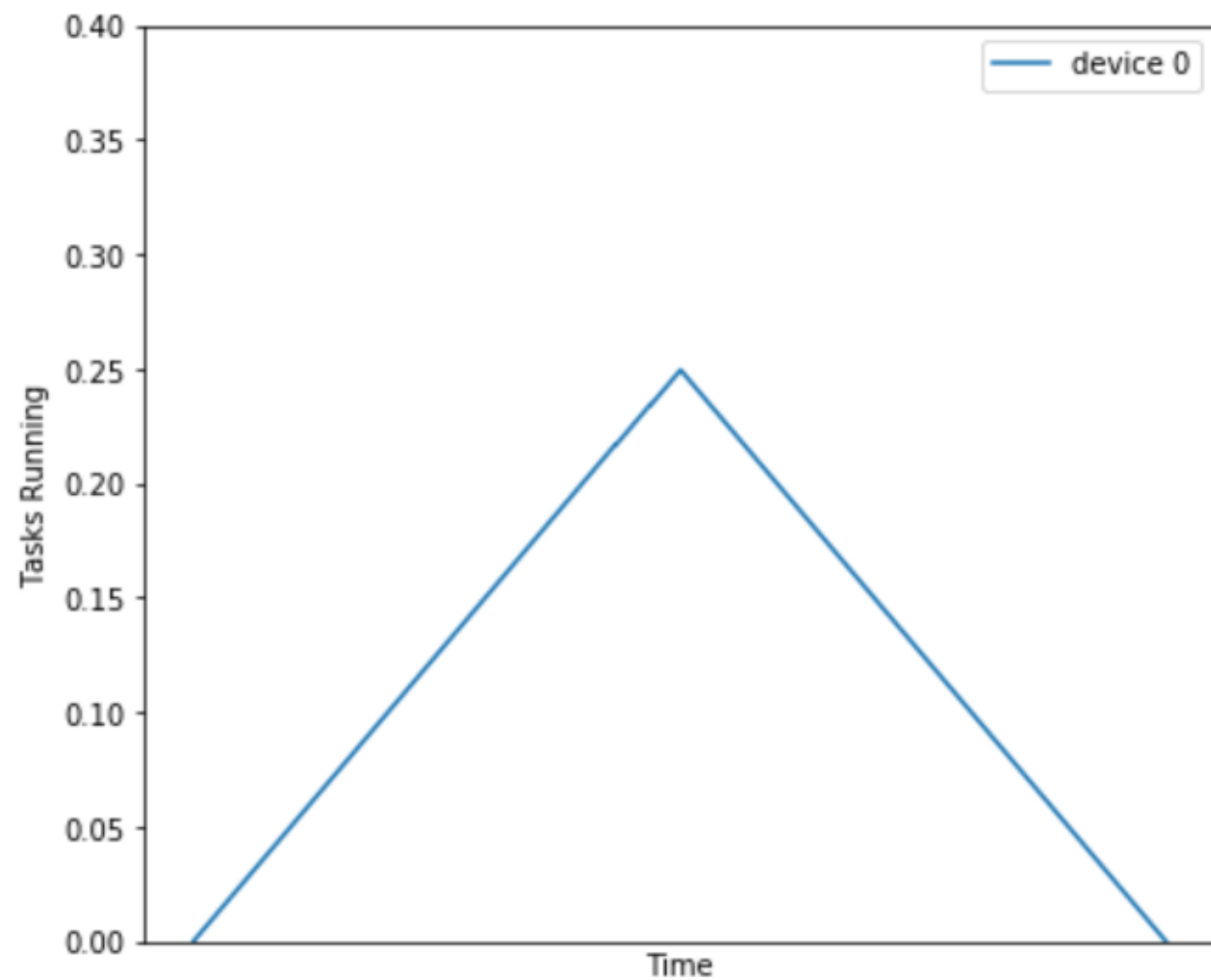
- 2 Intel Xeon CPUs
- 8 Nvidia Pascal GPUS
- 256 GB memory + 16 GB per GPU memory

## MCL Scheduler Comparison





# Round Robin vs. First Fit



The load that is just on GPU 0 in the First Fit scheduler is now balanced

## Another Problem

- Breadth First Search
  - Graph is represented in two arrays: one storing an adjacency list, one storing the offsets of each vertex – read only data structure
  - On each iteration  $i$ :
    - ✓ The frontier has an array of vertices that distance  $i$  from the source
    - ✓ The cost array has best known cost for each vertex
    - ✓ “Explore” indices that are  $i+1$  away from the source using the frontier array
- Each iteration is a new task reusing the same data from the previous tasks
- Memory transfers dominate compared to computation

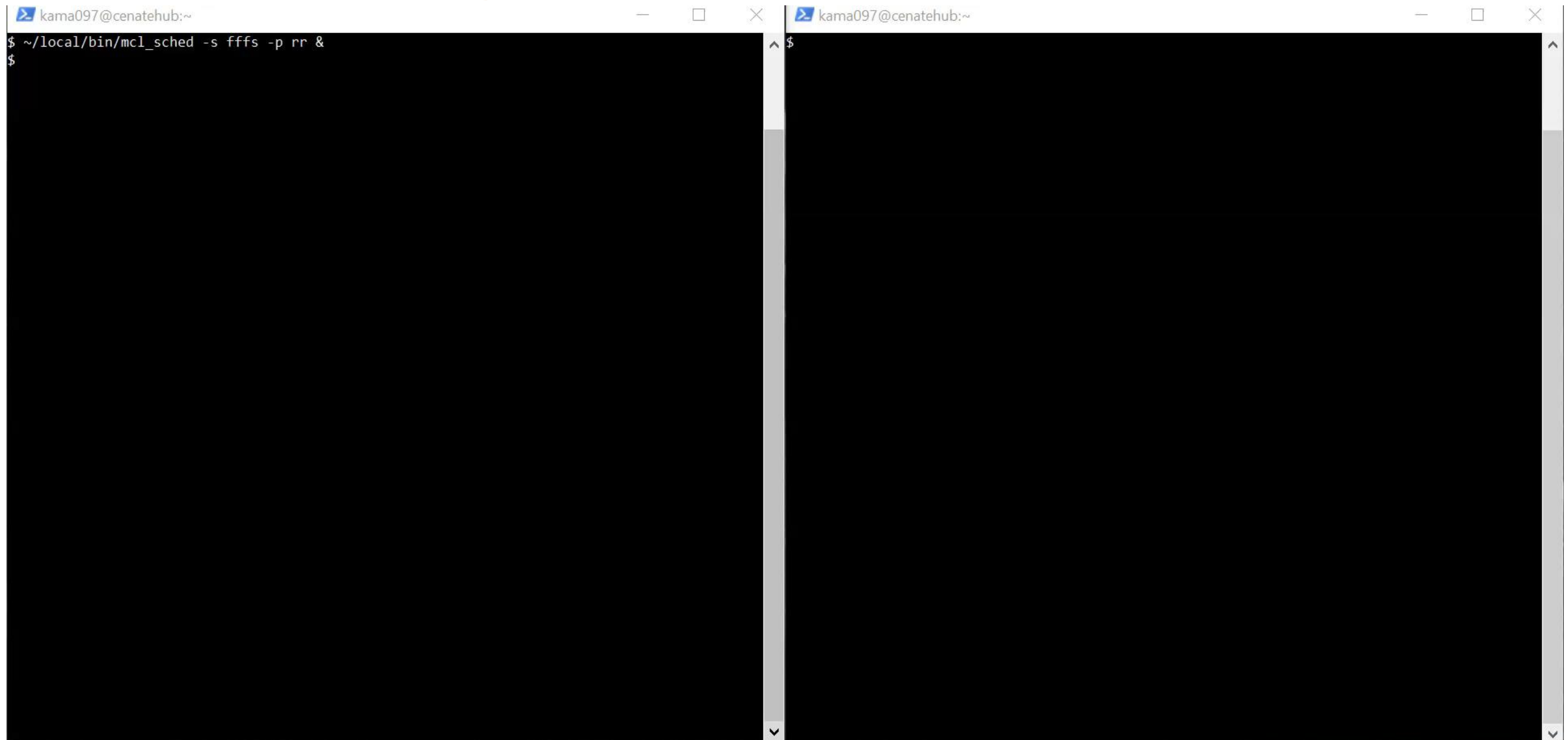
# Breadth First Search

- Breadth First Search Code per iterations:

```
mcl_handle* hdl = mcl_task_create();  
mcl_task_set_kernel(hdl, kernel_path, "BFS", 8, "", 0x0);  
mcl_task_set_arg(hdl, 0, (void*)frontier, ..., flags);  
mcl_task_set_arg(hdl, 1, (void*)edge_offsets, ..., flags);  
mcl_task_set_arg(hdl, 2, (void*)edge_list, ..., flags);  
...  
mcl_task_set_arg(hdl, 6, (void*)&iters, sizeof(uint32_t), MCL_ARG_SCALAR);  
mcl_exec(hdl, pes, lsize, MCL_TASK_ANY);
```



# Running Breadth First Search

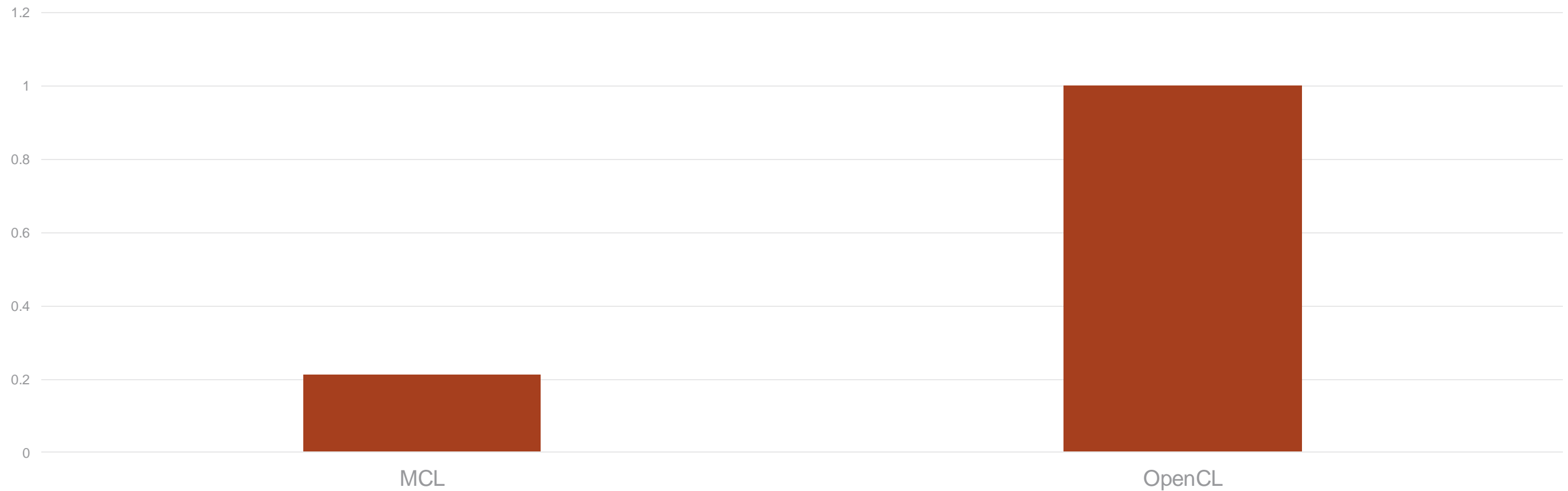


```
kama097@cenatehub:~  
$ ~/local/bin/mcl_sched -s fffs -p rr &  
$  
$
```

```
kama097@cenatehub:~  
$
```

# Breadth Fist Scheduler v1

MCL Performance vs. OpenCL



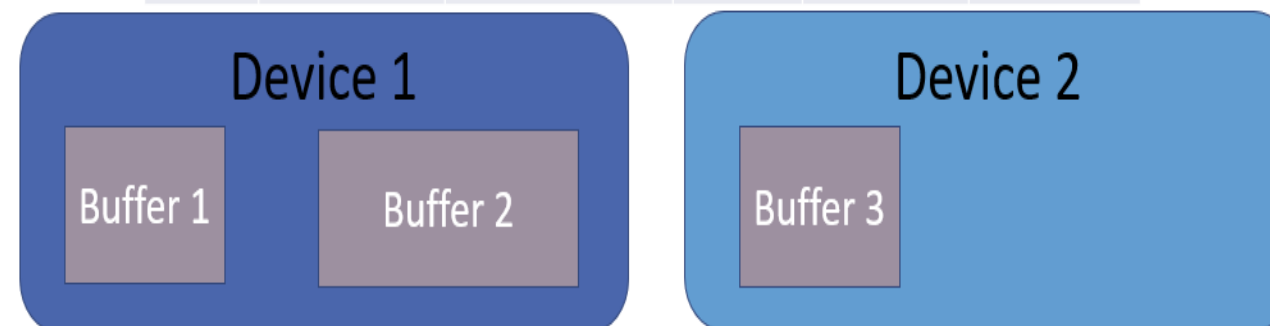
## Why so slow?

- Graph is naively transferred for each iteration

# MCL Resident Memory Module

- Allows persistent data to remain in device memory across tasks
- Orchestrates 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
...	...	...	...	...	...



MCL\_ARG\_RESIDENT

MCL\_ARG\_INVALID

MCL\_ARG\_DONE

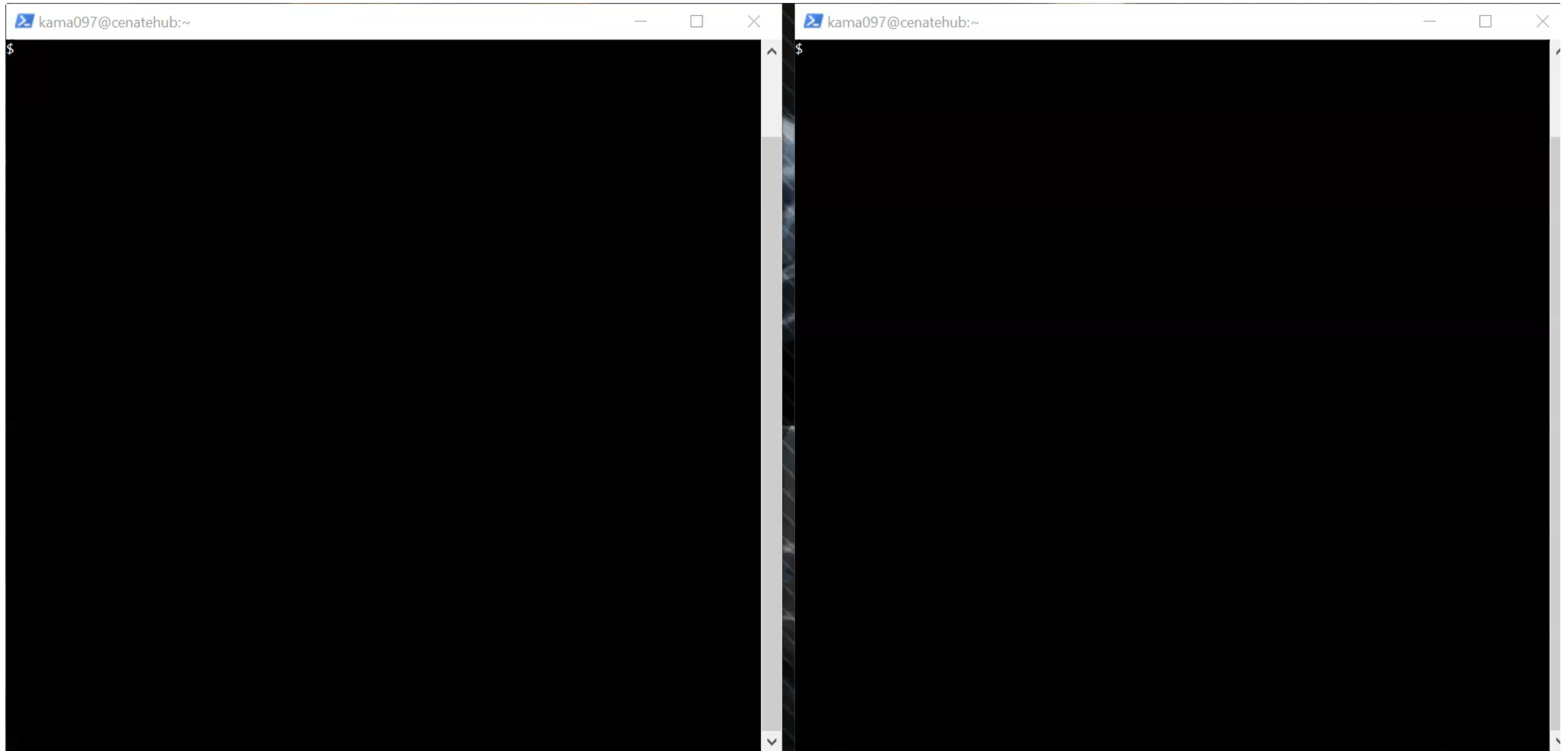


# Code Modification Demo

- Show (live?) modifications from non-resident memory to resident memory.
- If you are following along in the code, this is the difference between BFS.cpp and BFS-modified.cpp

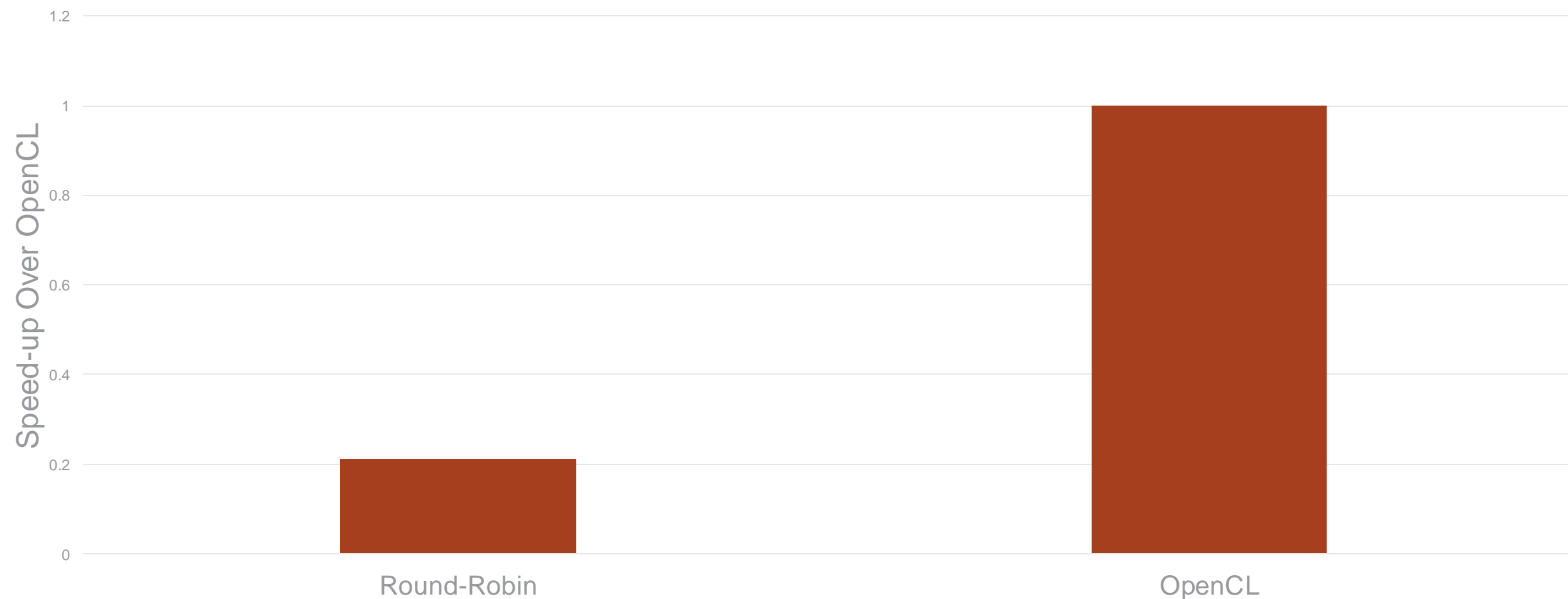
The image shows a Visual Studio Code interface with a dark theme. On the left, the Explorer sidebar displays a file tree for a project named 'ROOT [CONTAINER SHOC...]', listing various source files and headers. The main Editor window is open to 'BFS.cpp', showing C++ code for a Breadth-First Search (BFS) algorithm. The code includes comments in Chinese and uses macros for error checking and task setting. The status bar at the bottom indicates the current configuration, including 'multi-gpu\*', 'CMake: [Debug]: Ready', and the active language 'C++'.

# Running The New Code



# Breadth First Search v2

SHOC BFS Benchmark



## Still so slow?

- Different frontiers must be transferred from other devices

# Locality Aware Scheduler: Delay Scheduling

- Scheduler needs to choose the optimal device based off where resident memory is located
- Delay Scheduling:
  - 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
  - “Local data” is required data that is currently on a device

---

**Algorithm 1** Device Local Data

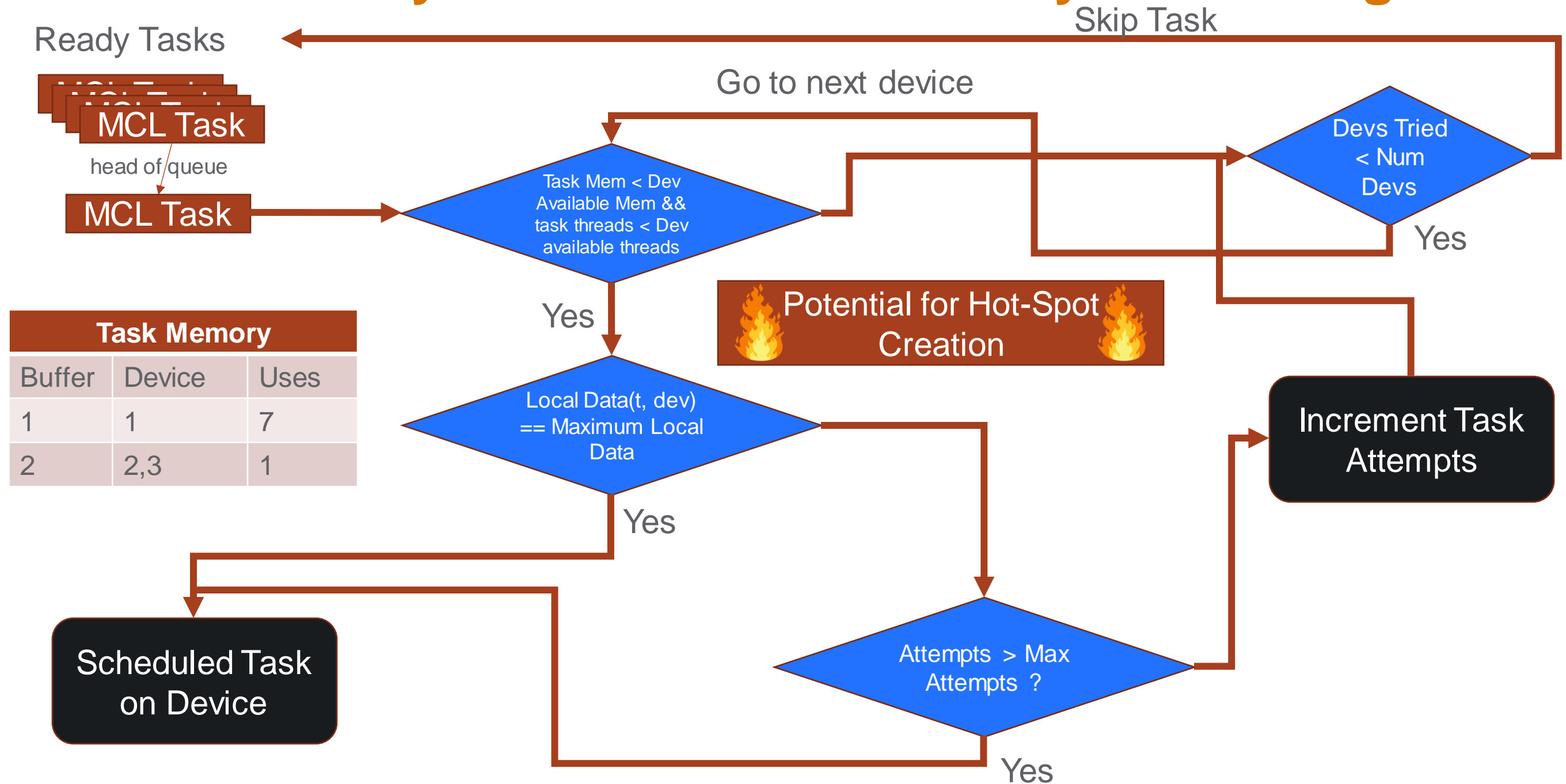
---

```
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
```

---

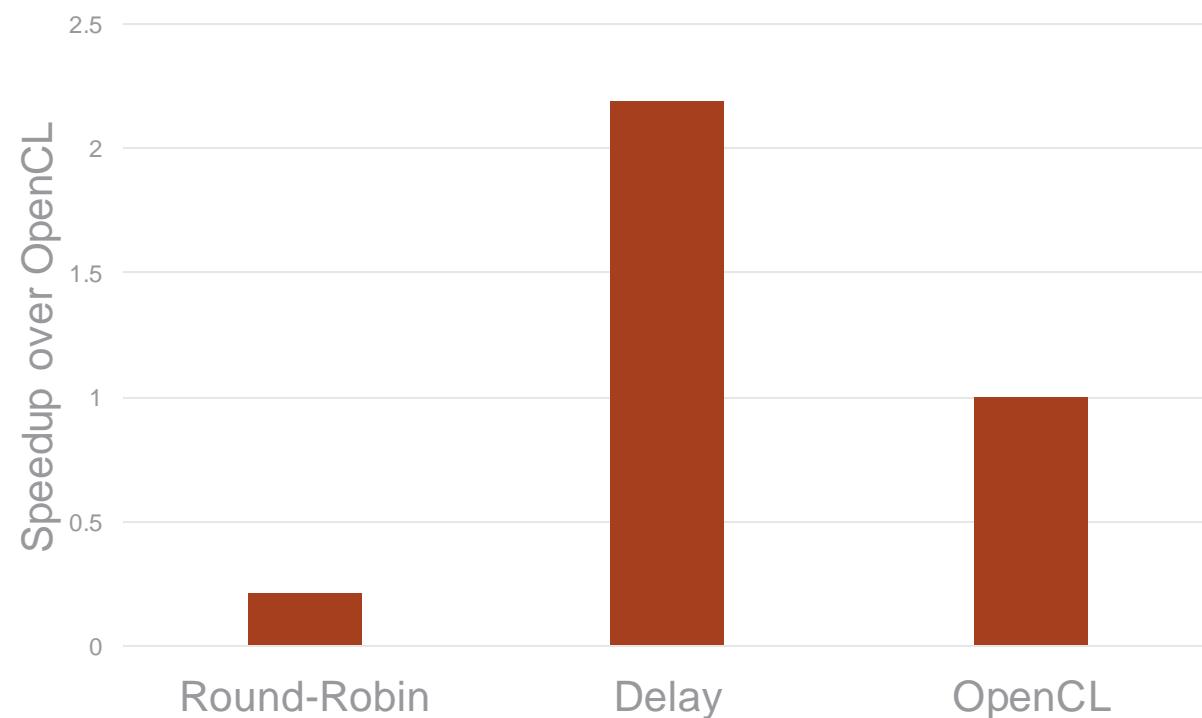


# Locality Aware Scheduler: Delay Scheduling

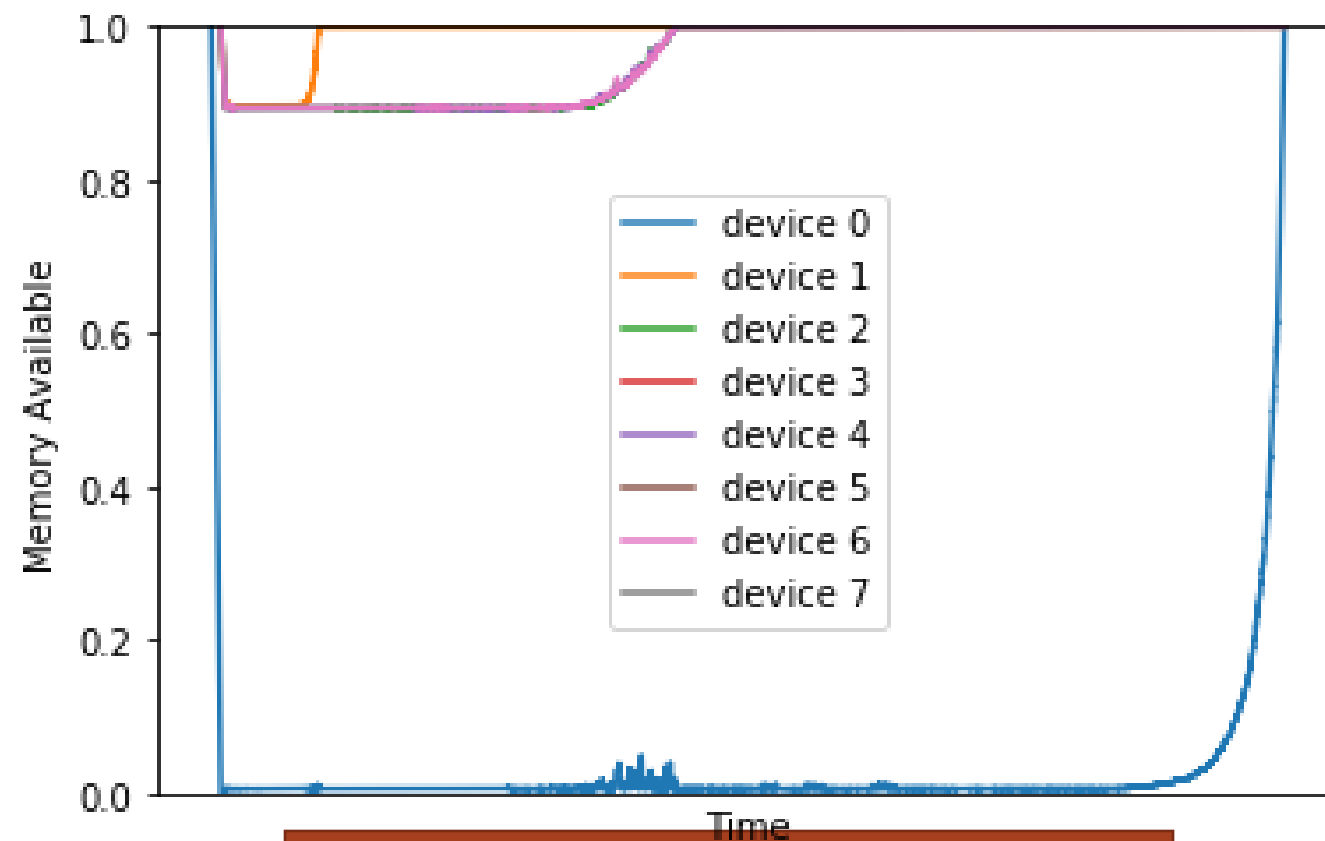


# Delay Scheduler Results

SHOC BFS Benchmark



Performance Improves! But can we do better...?



MCL scheduler trace reveals that one device is heavily used

# Hybrid Scheduler

- Needed to balance locality concerns against system utilization=
- Detects popular pieces of data to create replicas – done by changing how we calculate local data
- Hyperparameters are controlled with environment variables:  
MCL\_SCHED\_MAX\_ATTEMPTS  
and  
MCL\_SCHED\_COPY\_FACTOR

---

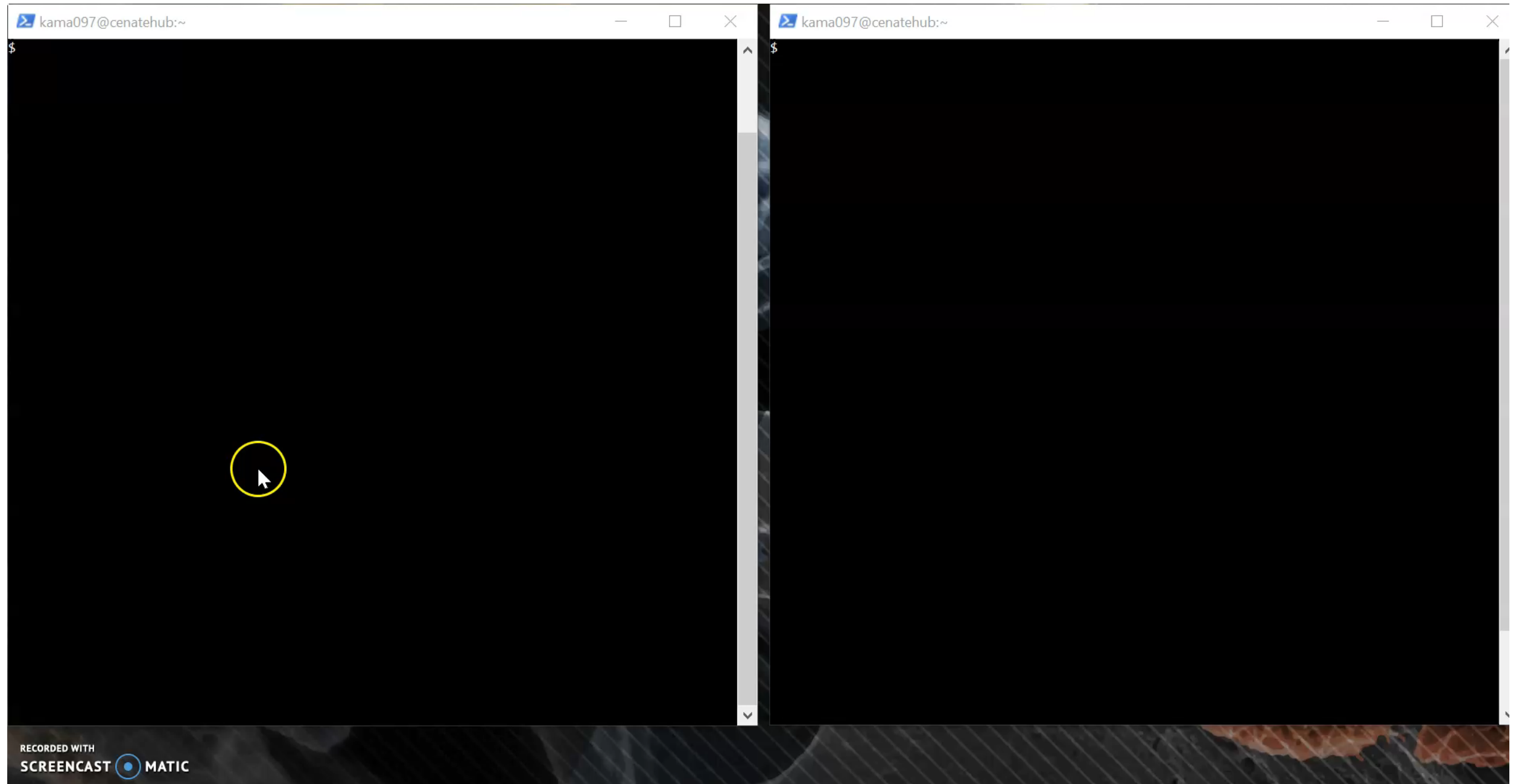
## Algorithm 3 Device Local Data With Copy Factor

---

```
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
```

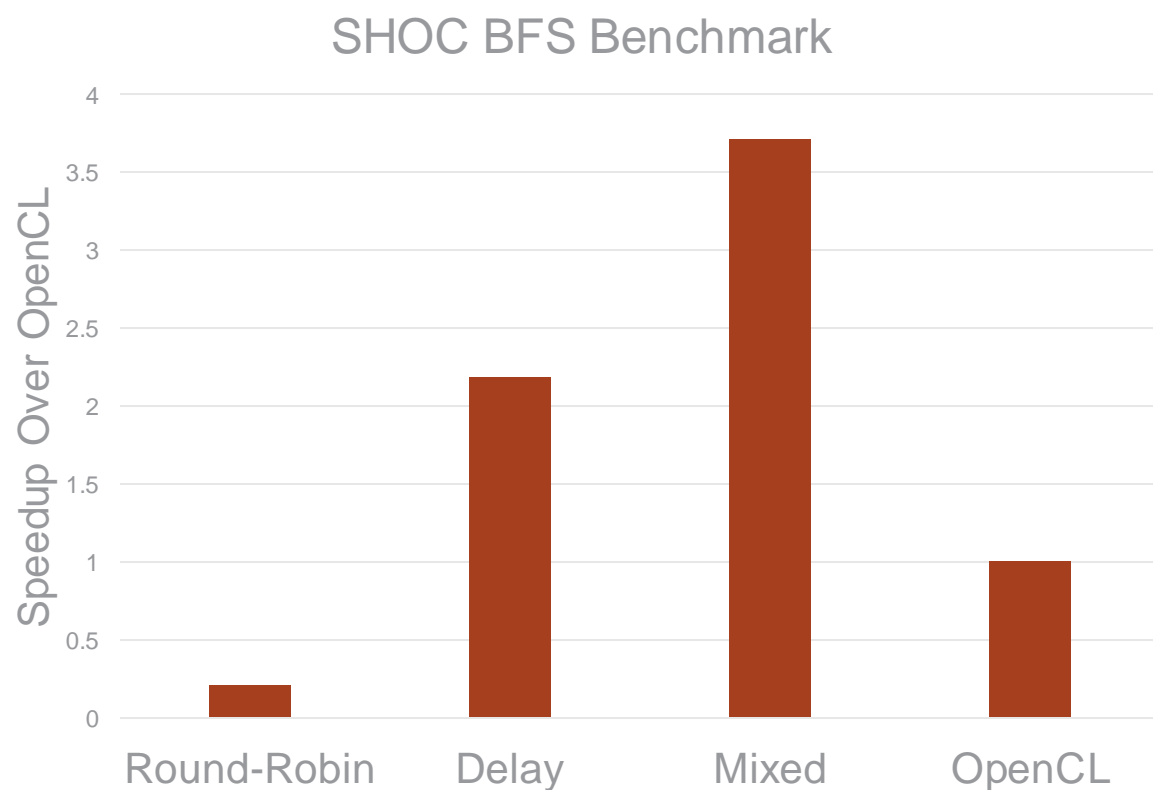
---

# BFS – Final Version Demo

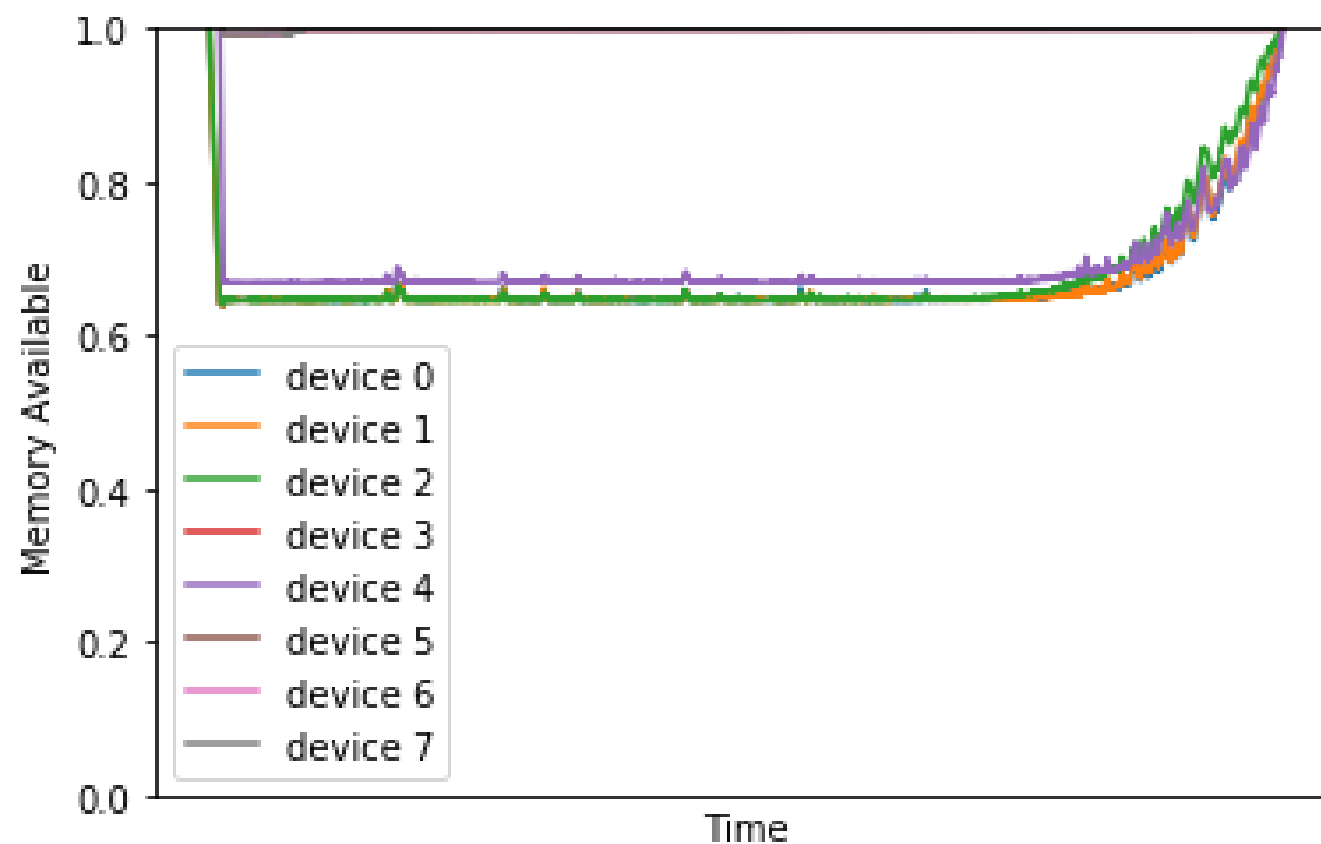




# Breadth First Search Results

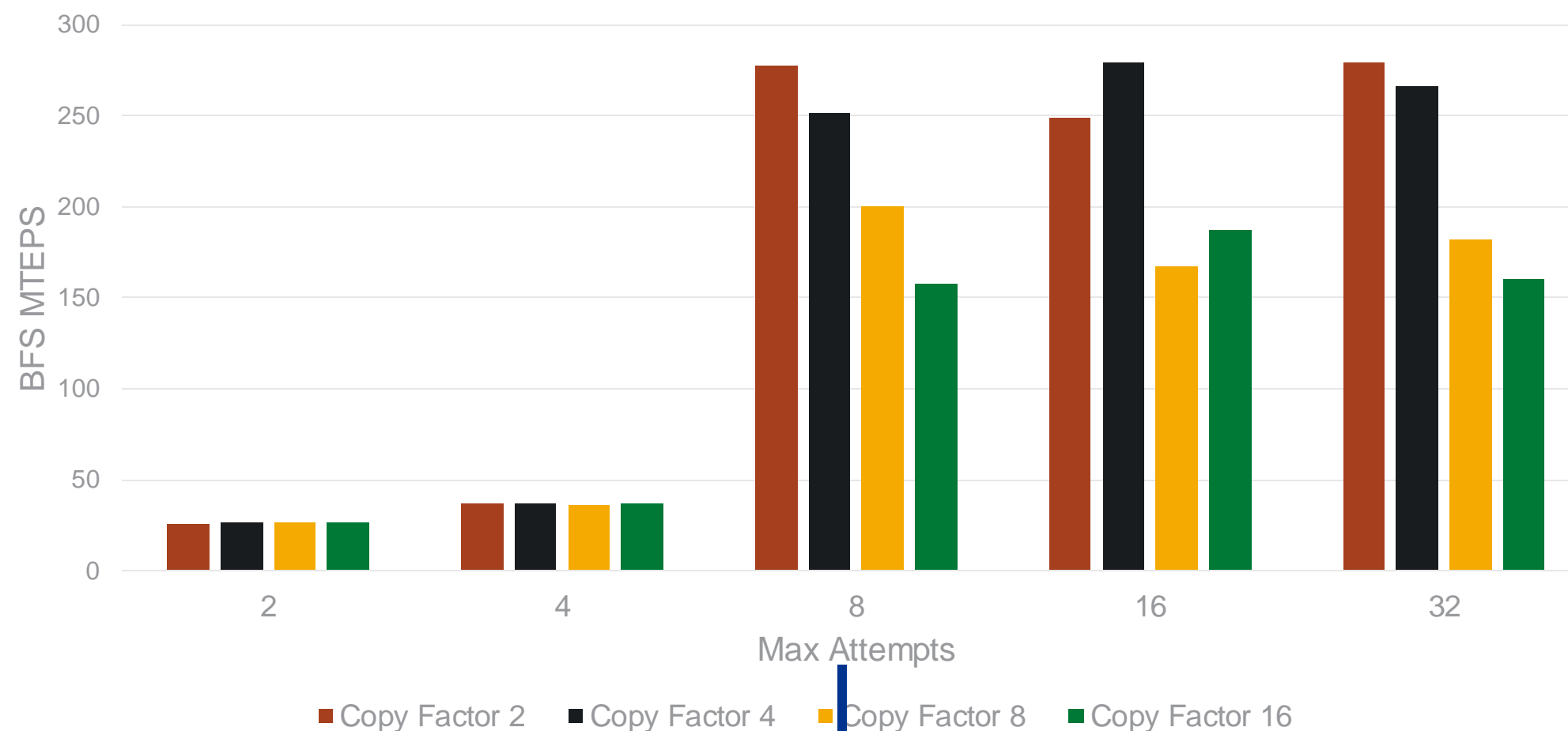


Best performance!



MCL scheduler trace reveals  
balance among first 4 devices

# Effect of Hyperparameters on Performance



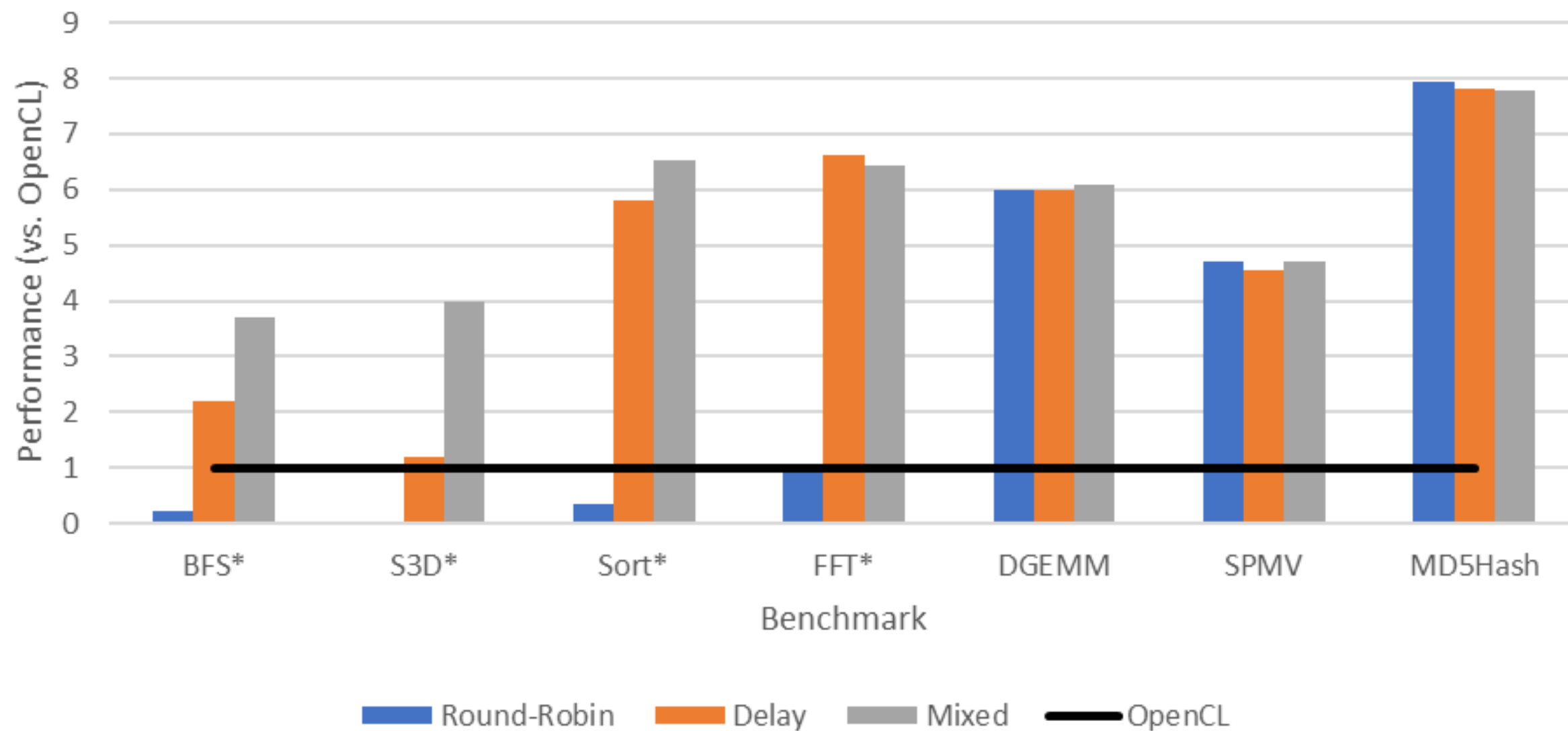
BFS  
Benchmark  
- 1,000,000  
vertices  
- 4096 Tasks

When Max Attempts <  
Num Devices – low  
performance

Num  
GPUs

If Max Attempts > Num  
Devices – copy factor  
dominates performance

# Full Scheduler Comparison



# Writing Your Own Scheduler

- Certain applications require specific requirements from the scheduler
- `init(mcl_resource_t* r, int ndevs)` – initialize representation of resource and any other representation needed
- `find_resource(sched_req_t r)` – Find the device `r` should run on. Also, set `r->dev` to the assigned device
- `assign_resource (sched_req_t r)` – Allocate resources in resource to the scheduled device
- `assign_resource (sched_req_t r)` – Release resources in resource model from scheduled device



# Eviction Policy

- 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

# Upcoming Work – Multi-Application Scheduling

- Different applications needs to be run in a pipeline
  - A physics simulation -> a data analysis application to detect events
  - A MD-simulation guided by reinforcement learning
  - Combustion Simulation + Machine Learning
  - Etc.
- Currently in 2 ways:
  - Modify exiting application -> specific to each application, lots of extra work,
  - Leverage Files

# Shared Memory Design

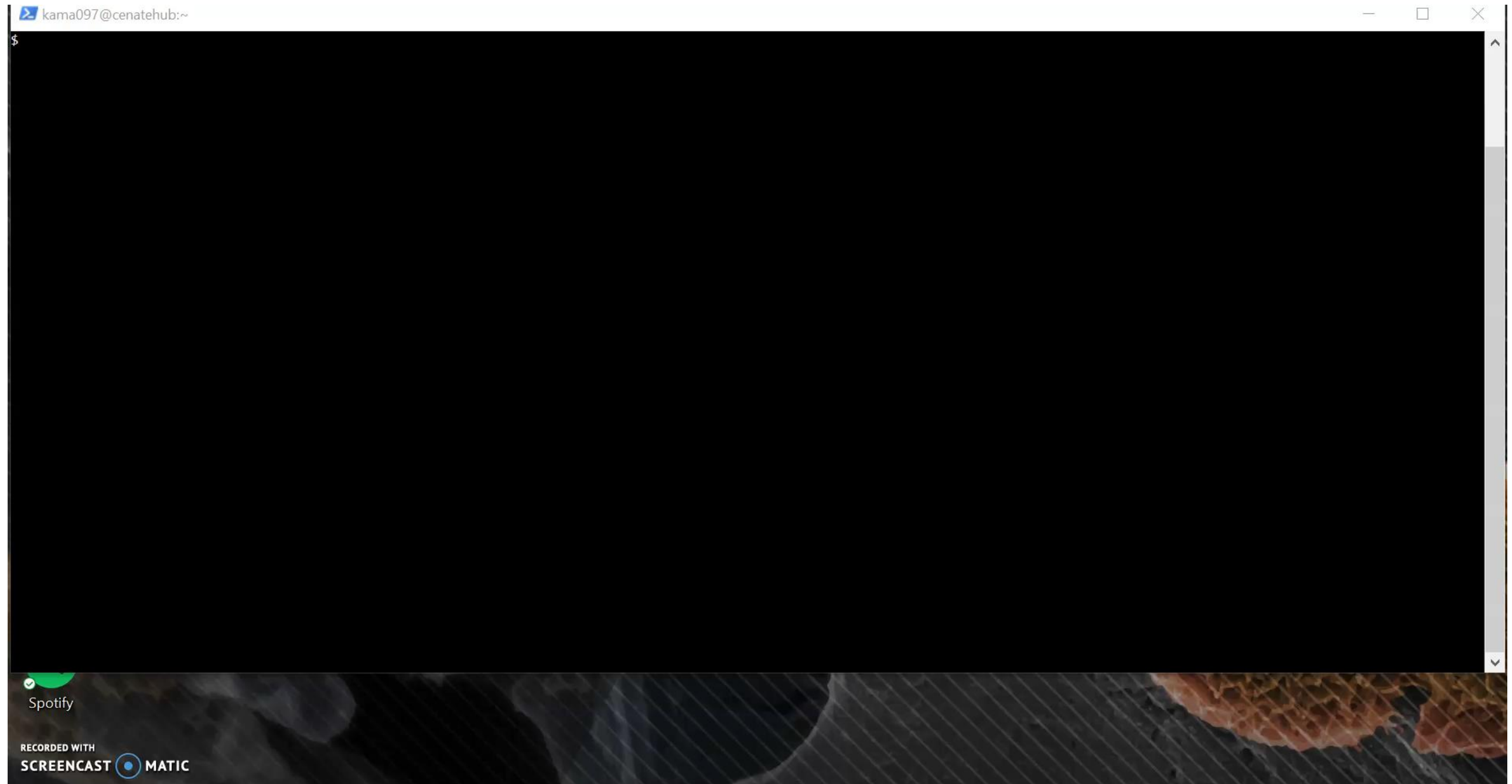
- Leverage POCL to create an additional OpenCL interface that allows buffers to be shared
- Scheduler is aware of where necessary data is even across applications
- The same data does not have to be transferred again in different applications
- Patterns Supported:
  - Scratch Pad
  - Producer-Consumer
  - Circular Buffer

# Pipelining Applications Code

- Experimental Application
  - Producer performs an arbitrary number of floating-point operations on a buffer
  - Consumer reads producers buffer and performs its own operations
- Comparisons:
  - MCL Shared memory
  - OpenCL – File
  - OpenCL – POSIX Shared Mem + Pipe (Statically Partitioned GPUs)



# Pipelined Applications (Example)



Thank you

