# AGENDA

What is GPU Direct?
CUDA Aware MPI
~~Advanced On Node Communication~~

# SUMMIT NODE OVERVIEW

# SUMMIT NODE

## (2) IBM POWER9 + (6) NVIDIA VOLTA V100



| 256 GB (DDR4) | | | | | 256 GB (DDR4) | | |
|---|---|---|---|---|---|---|---|

135 GB/s

CPU 0

| 0 (0-3) | 7 (28-31) | 14 (56-59) |
|---|---|---|
| 1 (4-7) | 8 (32-35) | 15 (60-63) |
| 2 (8-11) | 9 (36-39) | 16 (64-67) |
| 3 (12-15) | 10 (40-43) | 17 (68-71) |
| 4 (16-19) | 11 (44-47) | 18 (72-75) |
| 5 (20-23) | 12 (48-51) | 19 (76-79) |
| 6 (24-27) | 13 (52-55) | 20 (80-83) |

64 GB/s

135 GB/s

CPU 1

| 22 (88-91) | 29 (116-119) | 36 (144-147) |
|---|---|---|
| 23 (92-95) | 30 (120-123) | 37 (148-151) |
| 24 (96-99) | 31 (124-127) | 38 (152-155) |
| 25 (100-103) | 32 (128-131) | 39 (156-159) |
| 26 (104-107) | 33 (132-135) | 40 (160-163) |
| 27 (108-111) | 34 (136-139) | 41 (164-167) |
| 28 (112-115) | 35 (140-143) | 42 (168-171) |

GPU 0   GPU 1   GPU 2

16 GB (HBM2)   16 GB (HBM2)   16 GB (HBM2)

GPU 3   GPU 4   GPU 5

16 GB (HBM2)   16 GB (HBM2)   16 GB (HBM2)

NVLink2 ⟷ (50 GB/s)   ↕ (900 GB/s)

NVIDIA.

# UNDER THE HOOD

Summit has fat nodes!

Many connections

Many devices

Many stacks



| | | |
|---|---|---|
| TF | 42 TF (6x7 TF) | |
| HBM | 96 GB (6x16 GB) | |
| DRAM | 512 GB (2x16x16 GB) | |
| NET | 25 GB/s (2x12.5 GB/s) | |
| MMsg/s | 83 | |

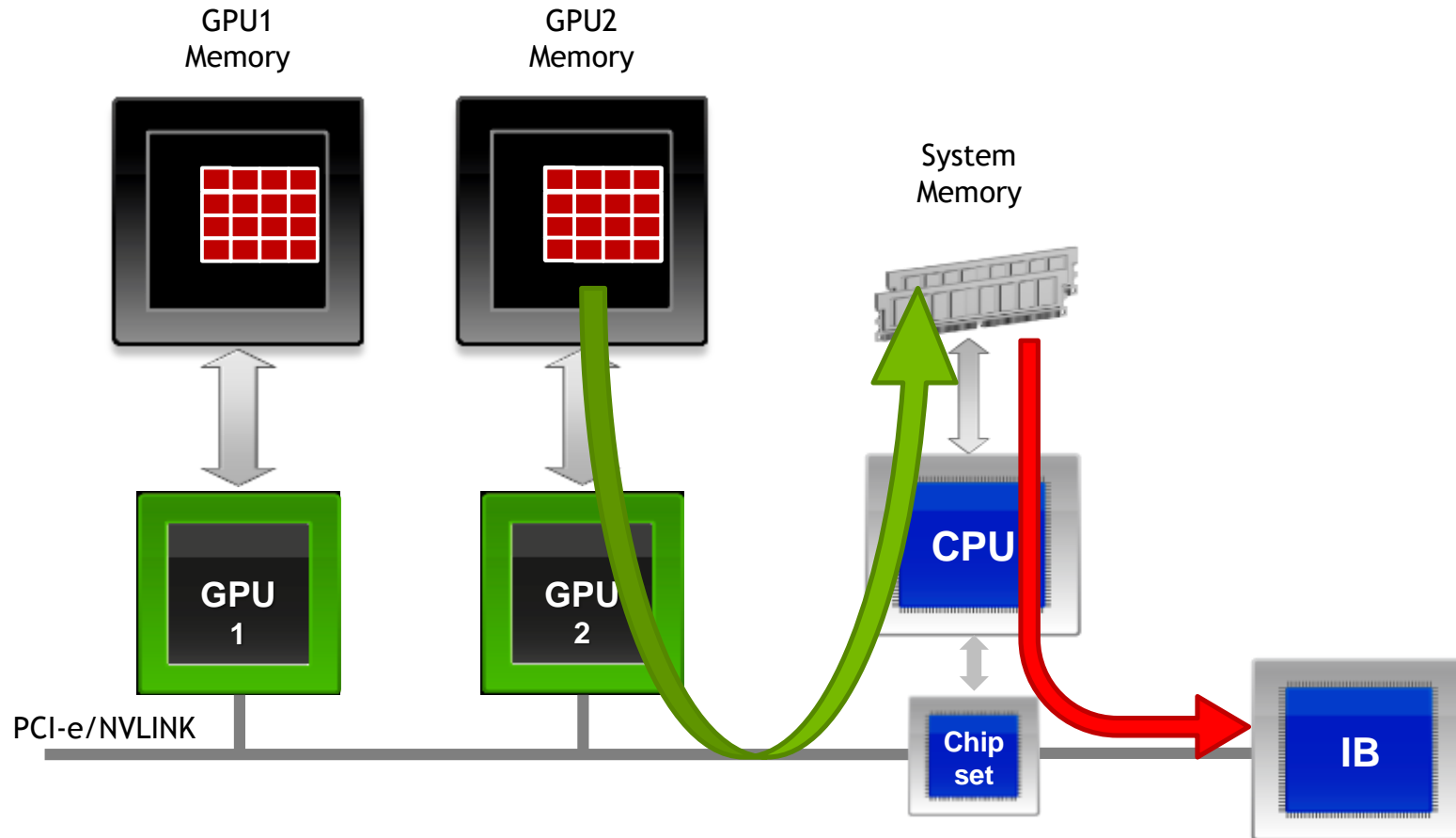- HBM/DRAM Bus (aggregate B/W)
- NVLINK
- X-Bus (SMP)
- PCIe Gen4
- EDR IB

HBM & DRAM speeds are aggregate (Read+Write).
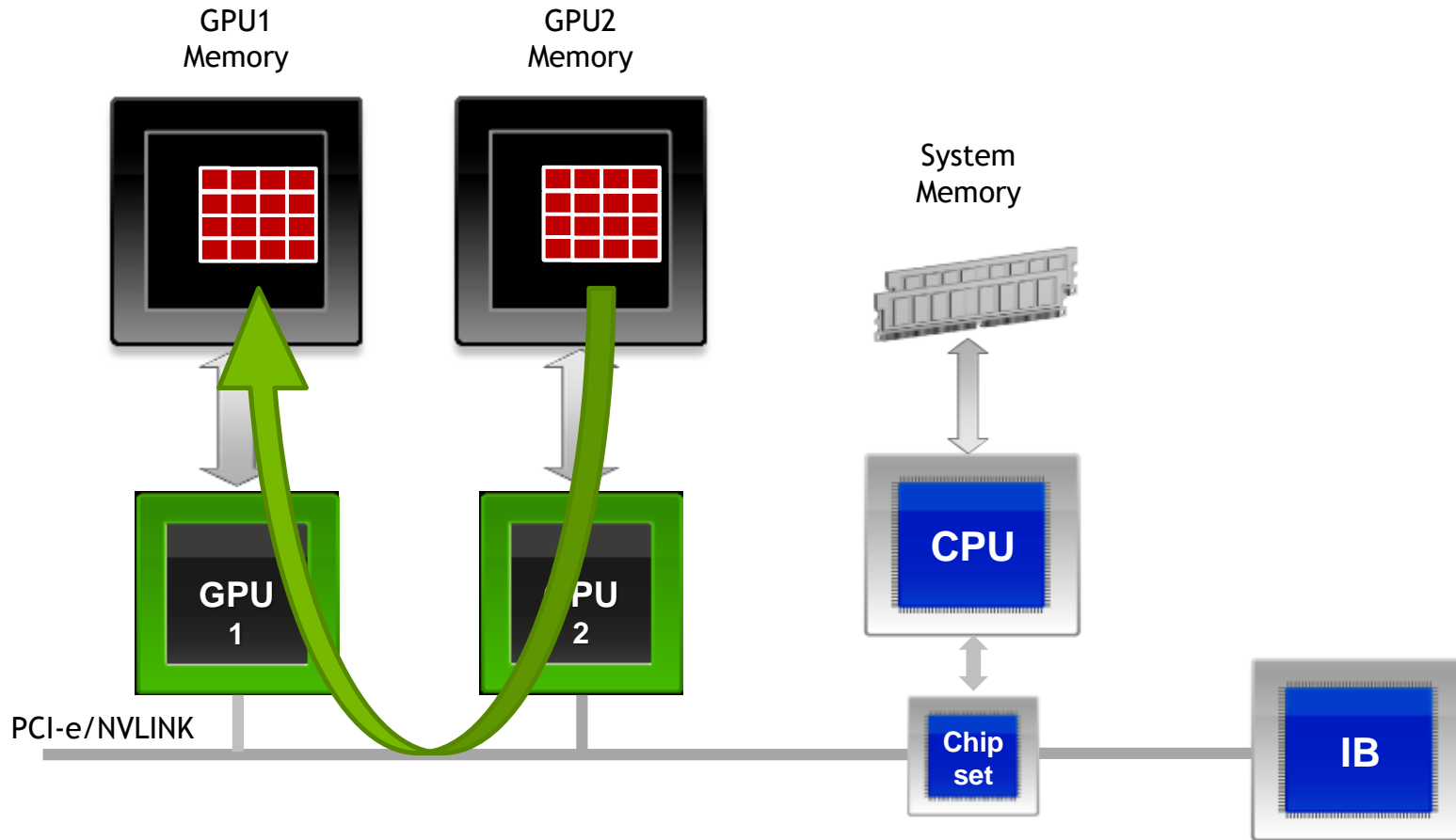All other speeds (X-Bus, NVLink, PCIe, IB) are bi-directional.

# GPUDIRECT

# NVIDIA GPUDIRECT™

## Accelerated Communication with Network & Storage Devices
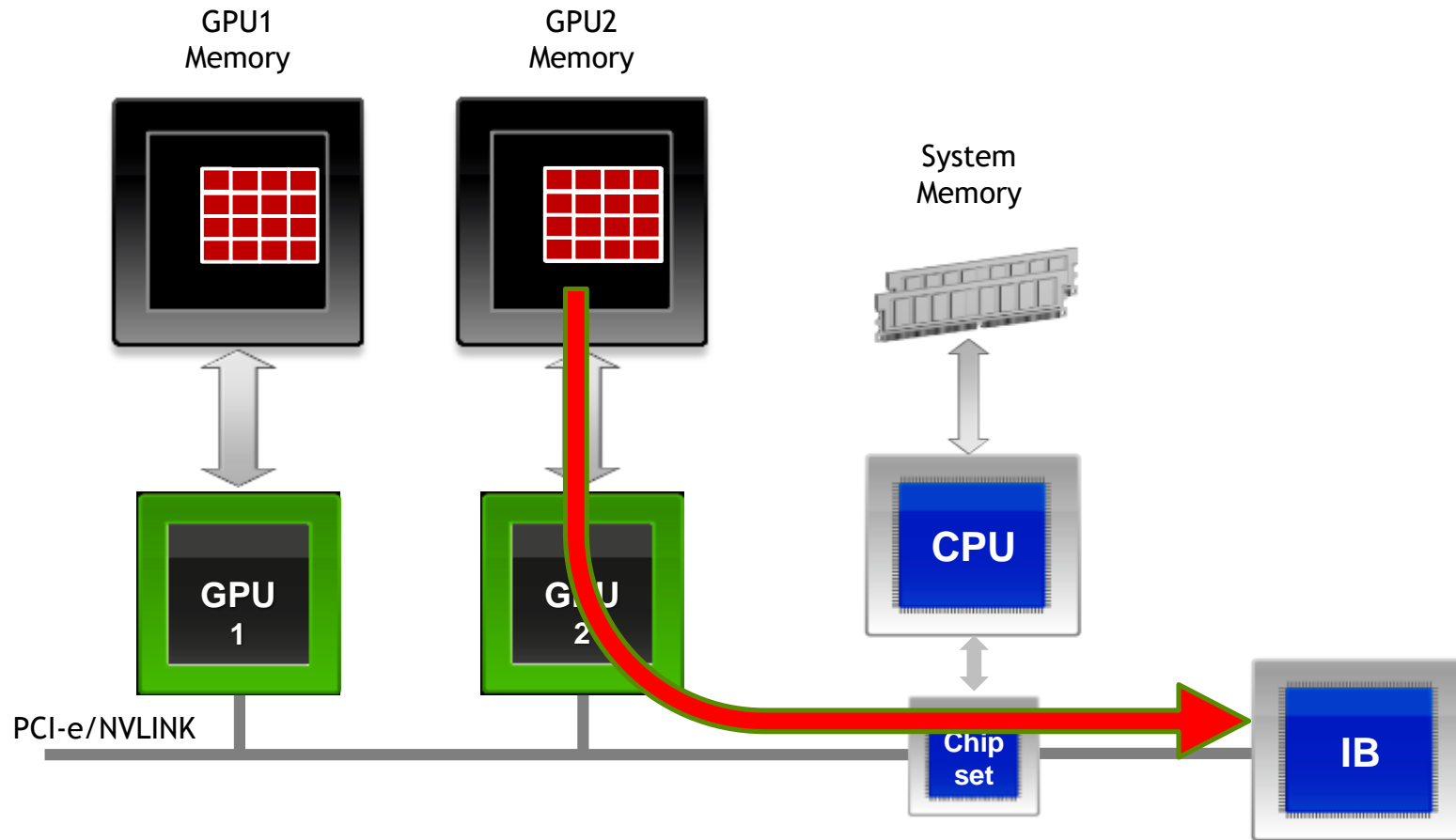
# NVIDIA GPUDIRECT™

## Peer to Peer Transfers

GPU1
Memory

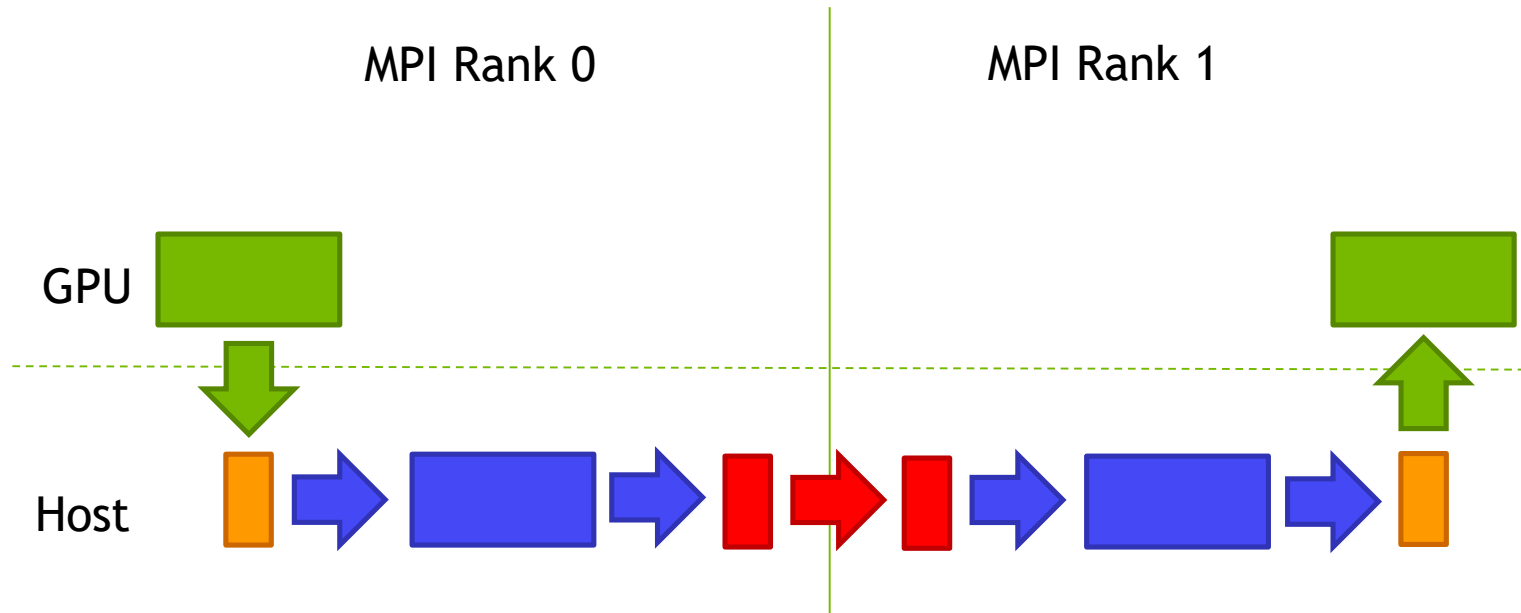GPU2
Memory

System
Memory

CPU

GPU
1

PU
2

Chip
set

IB

PCI-e/NVLINK

# CUDA AWARE MPI FOR ON AND OFF NODE TRANSFERS

# REGULAR MPI GPU TO REMOTE GPU

MPI Rank 0

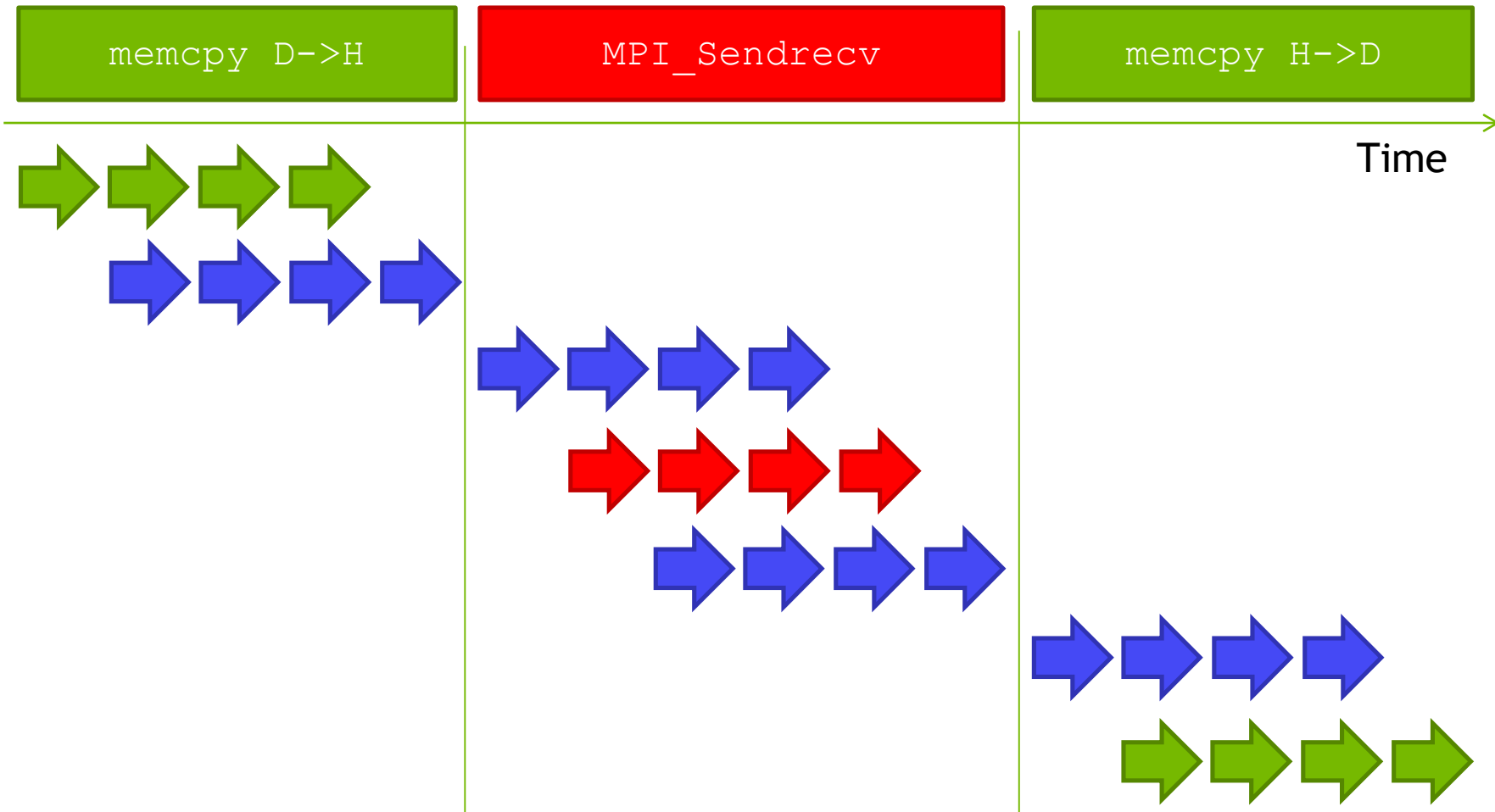MPI Rank 1

GPU

Host

```
cudaMemcpy(s_buf_h,s_buf_d,size,cudaMemcpyDeviceToHost);
MPI_Send(s_buf_h,size,MPI_CHAR,1,tag,MPI_COMM_WORLD);

MPI_Recv(r_buf_h,size,MPI_CHAR,0,tag,MPI_COMM_WORLD,&stat);
cudaMemcpy(r_buf_d,r_buf_h,size,cudaMemcpyHostToDevice);
```
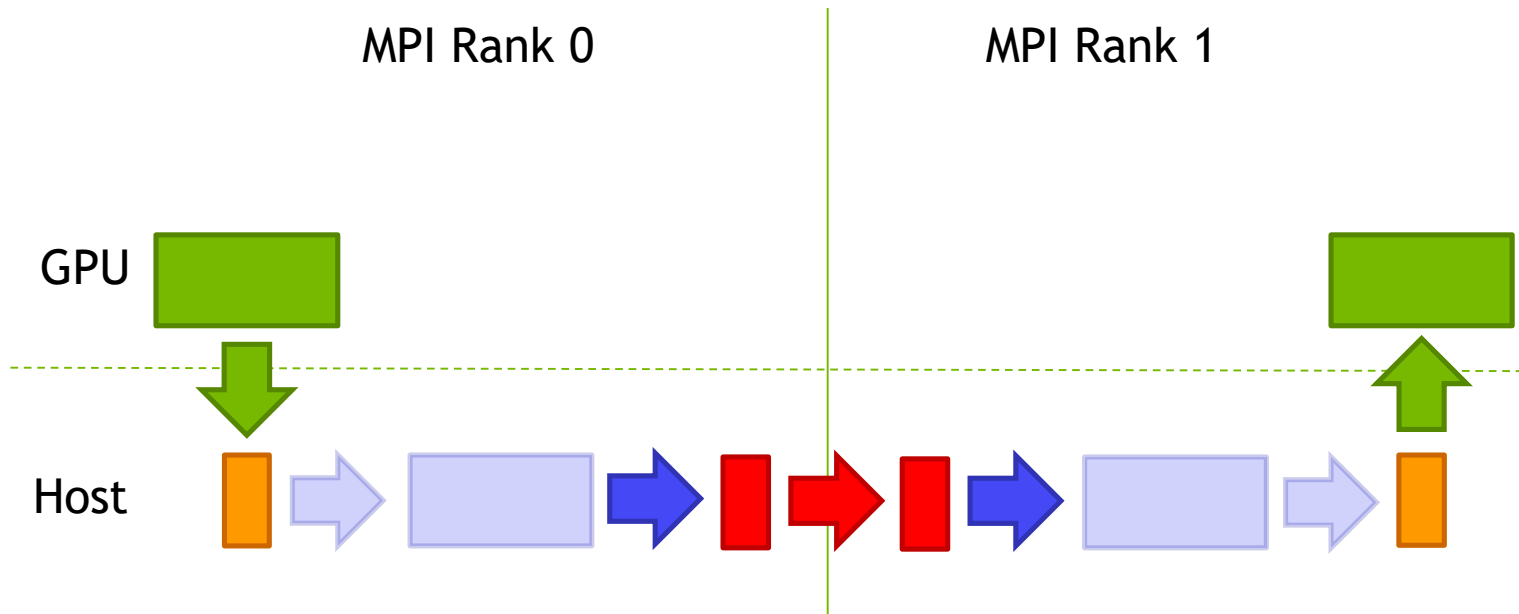
# REGULAR MPI GPU TO REMOTE GPU

| memcpy D->H | MPI_Sendrecv | memcpy H->D |

Time

NVIDIA.

# MPI GPU TO REMOTE GPU

## without GPUDirect

MPI Rank 0                    MPI Rank 1

GPU

Host

```
MPI_Send(s_buf_d,size,MPI_CHAR,1,tag,MPI_COMM_WORLD);

MPI_Recv(r_buf_d,size,MPI_CHAR,0,tag,MPI_COMM_WORLD,&stat);
```

NVIDIA.

# MPI GPU TO REMOTE GPU

## without GPUDirect

MPI Rank 0                    MPI Rank 1

GPU

Host

```
#pragma acc host_data use_device (s_buf, r_buf)
MPI_Send(s_buf,size,MPI_CHAR,1,tag,MPI_COMM_WORLD);

MPI_Recv(r_buf,size,MPI_CHAR,0,tag,MPI_COMM_WORLD,&stat);
```

# MPI GPU TO REMOTE GPU
## without GPUDirect

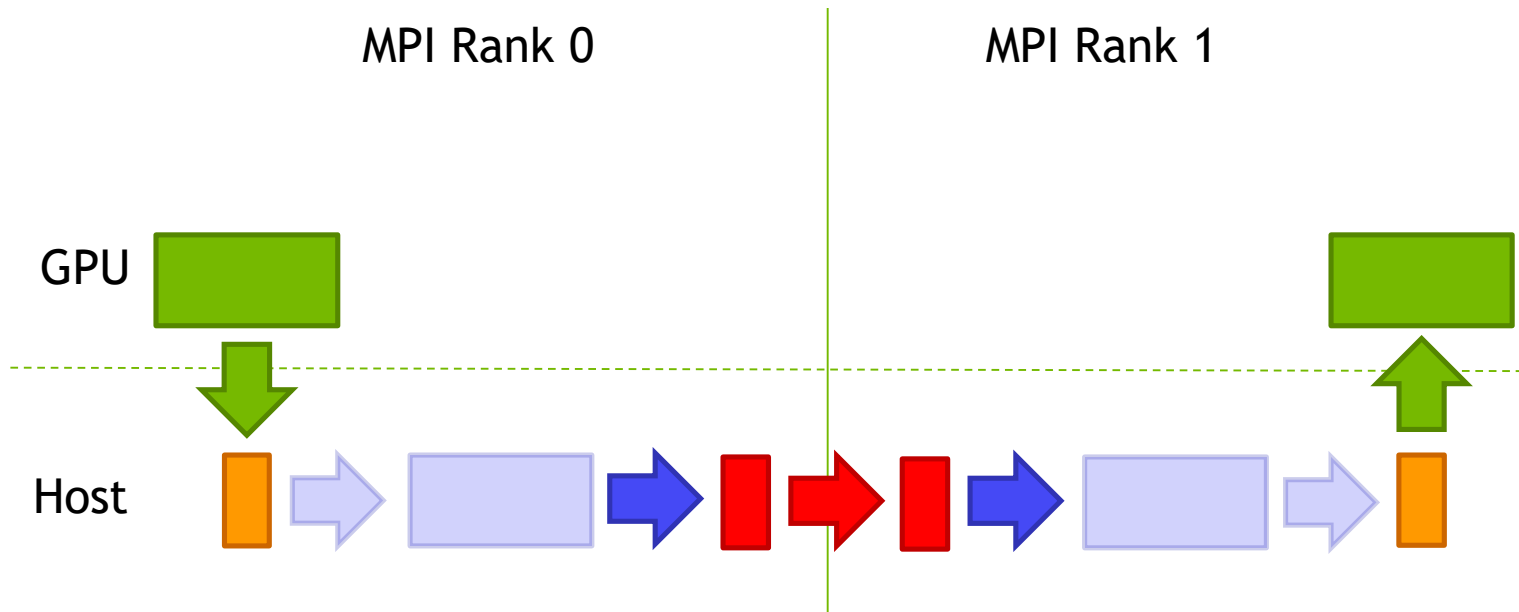# MPI GPU TO REMOTE GPU

## Support for RDMA

MPI Rank 0

MPI Rank 1

GPU

Host

```
MPI_Send(s_buf_d,size,MPI_CHAR,1,tag,MPI_COMM_WORLD);

MPI_Recv(r_buf_d,size,MPI_CHAR,0,tag,MPI_COMM_WORLD,&stat);
```

# MPI GPU TO REMOTE GPU

## Support for RDMA

MPI Rank 0

MPI Rank 1

GPU

Host
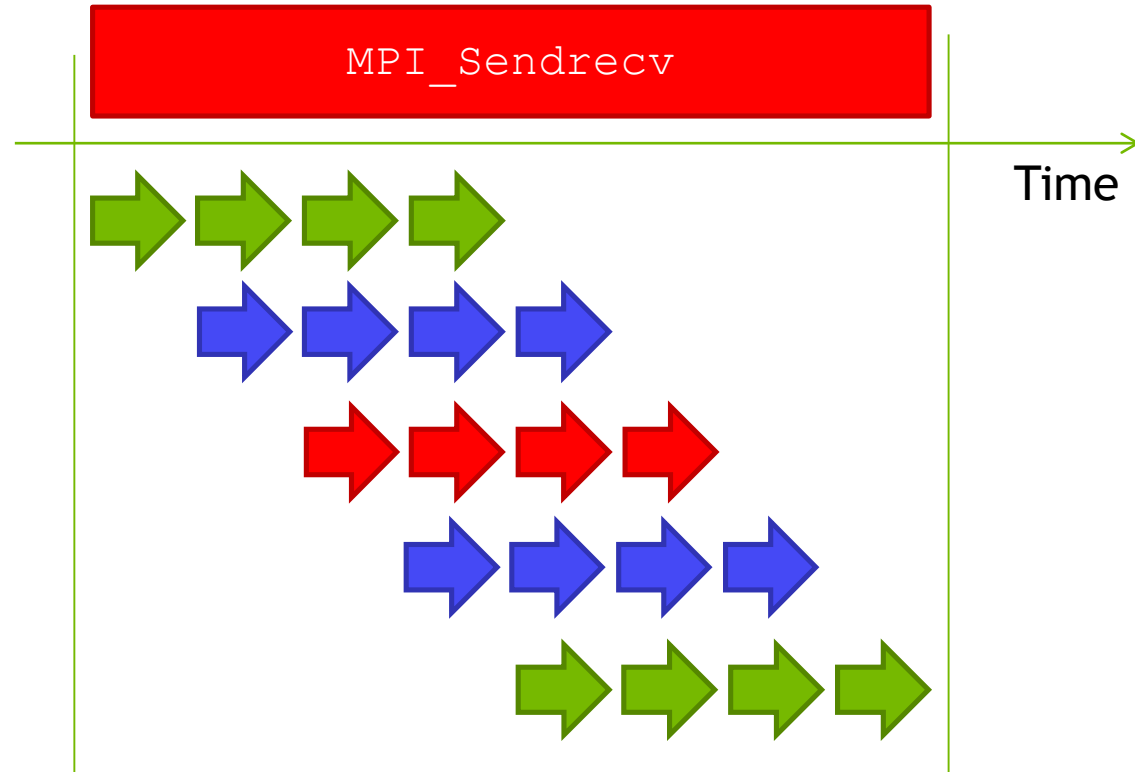
```
#pragma acc host_data use_device (s_buf, r_buf)
MPI_Send(s_buf,size,MPI_CHAR,1,tag,MPI_COMM_WORLD);

MPI_Recv(r_buf,size,MPI_CHAR,0,tag,MPI_COMM_WORLD,&stat);
```

17

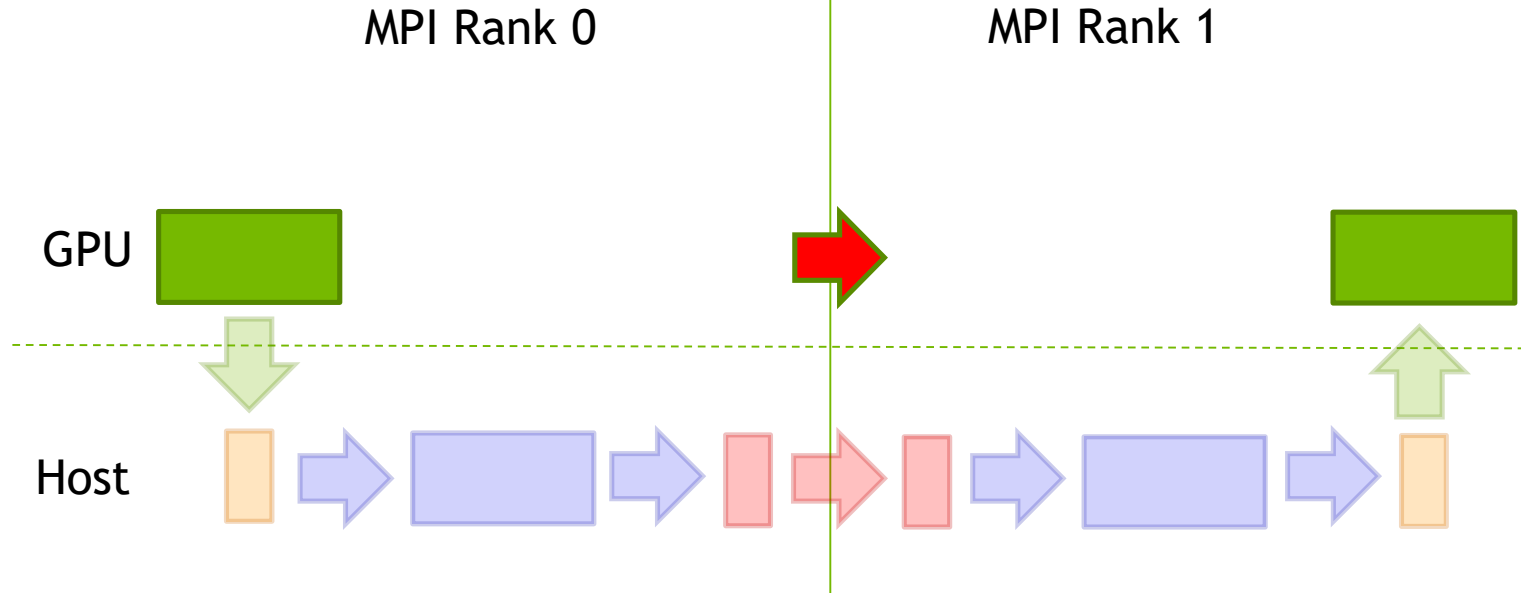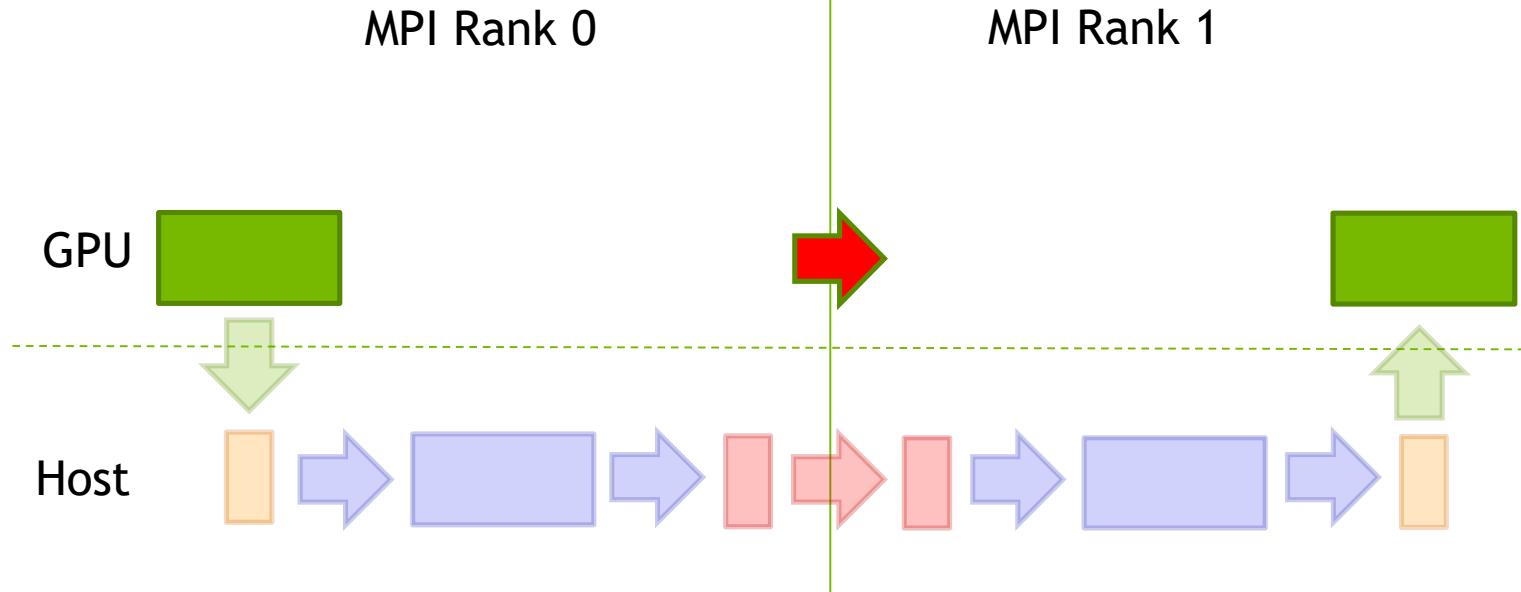# MPI GPU TO REMOTE GPU
## Support for RDMA

# JSRUN/SMPI GPU OPTIONS
## Running On Summit

To enable CUDA aware MPI, use **jsrun --smpiargs="-gpu"**

# KNOWN ISSUES ON SUMMIT
## Things to watch out for (as of January)

**Problems with Multiple resource sets per node:**

$> jsrun –g 1 –a 1 --smpiargs="-gpu" ….

[1]Error opening IPC Memhandle from peer:0, invalid argument

- One workaround: set PAMI_DISABLE_IPC=1

  - Expect poor performance, but a good functionality check

Will be resolved by software updates later this year

# PERFORMANT WORKAROUNDS
## Running On Summit

**Option 1: Run in one resource set and set GPU affinity in your code**

> **(do NOT restrict CUDA_VISIBLE_DEVICES, but you can permute it)**

**Option 2: Use a wrapper script**

- **Add "#BSUB –step_cgroup n" to your LSF options**

- **Run with `jsrun <your-jsrun-options> --smpiargs="-gpu" ./gpu_setter.sh <your app>`**

  - **(script on next slide)**

- **Will need to be careful about your CPU bindings!**

```bash
#! /bin/bash
# gpu_setter.sh
# Rudimentary GPU affinity setter for Summit
# >$ jsrun -rs_per_host 1 -gpu_per_rs 6 <task/cpu option> ./gpu_setter.sh <your app>

# This script assumes your code does not attempt to set its own
# GPU affinity (e.g. with cudaSetDevice). Using this affinity script
# with a code that does its own internal GPU selection probably won't work!

# Compute device number from OpenMPI local rank environment variable
# Keeping in mind Summit has 6 GPUs per node

mydevice=$((${OMPI_COMM_WORLD_LOCAL_RANK} % 6))

# CUDA_VISIBLE_DEVICES controls both what GPUs are visible to your process
# and the order they appear in. By putting "mydevice" first the in list, we
# make sure it shows up as device "0" to the process so it's automatically selected.
# The order of the other devices doesn't matter, only that all devices (0-5) are
# present.

CUDA_VISIBLE_DEVICES="${mydevice},0,1,2,3,4,5"

# Process with sed to remove the duplicate and reform the list, keeping the order we
set

CUDA_VISIBLE_DEVICES=$(sed -r ':a; s/\b([[:alnum:]]+)\b(.*)\b\1\b/\1\2/g; ta;
s/(,,)+/,/g; s/, *$//' <<< $CUDA_VISIBLE_DEVICES)

export CUDA_VISIBLE_DEVICES

# Launch the application we were given
exec "$@"
```

# GPU TO GPU COMMUNICATION

- CUDA aware MPI functionally portable

  - OpenACC/MP interoperable

  - Performance may vary between on/off node, socket, HW support for GPU Direct

  - Unified memory support varies between implementations, but it becoming common

- For more information on the following advanced on-node communication models, see the OLCF Training Archive or Summit Training Workshops

  - Single-process, multi-GPU

  - Multi-process, single-gpu