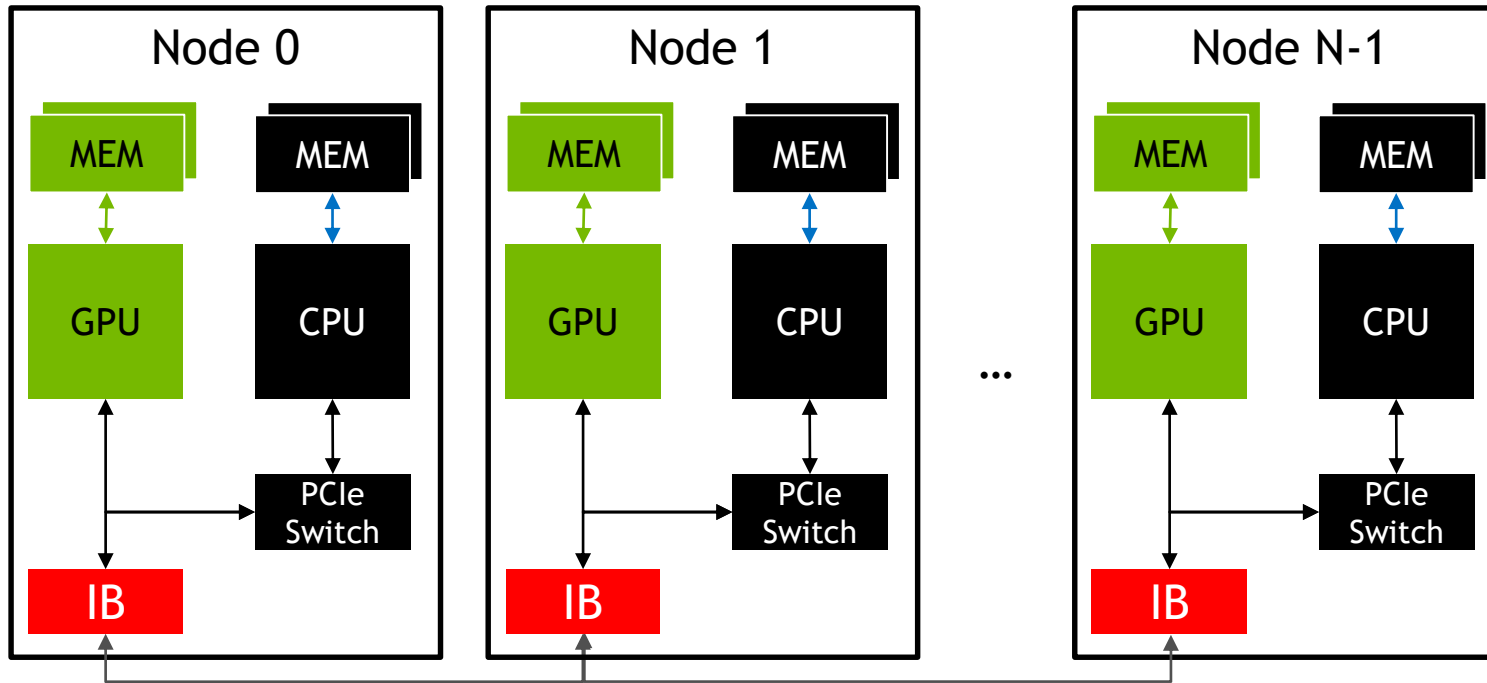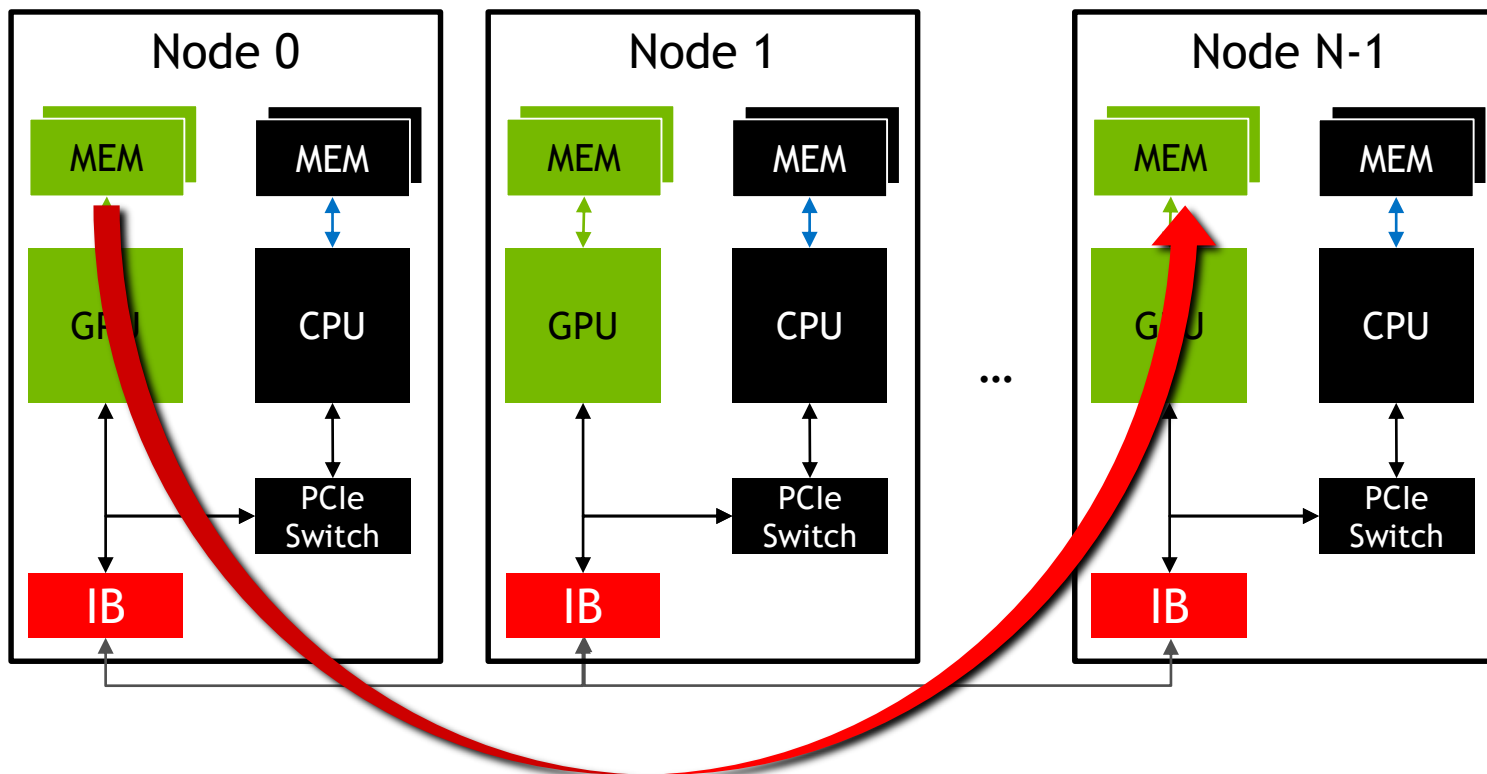# Multi-GPU Programming with MPI

Cheng Yi, Senior Solution Architect, NVIDIA
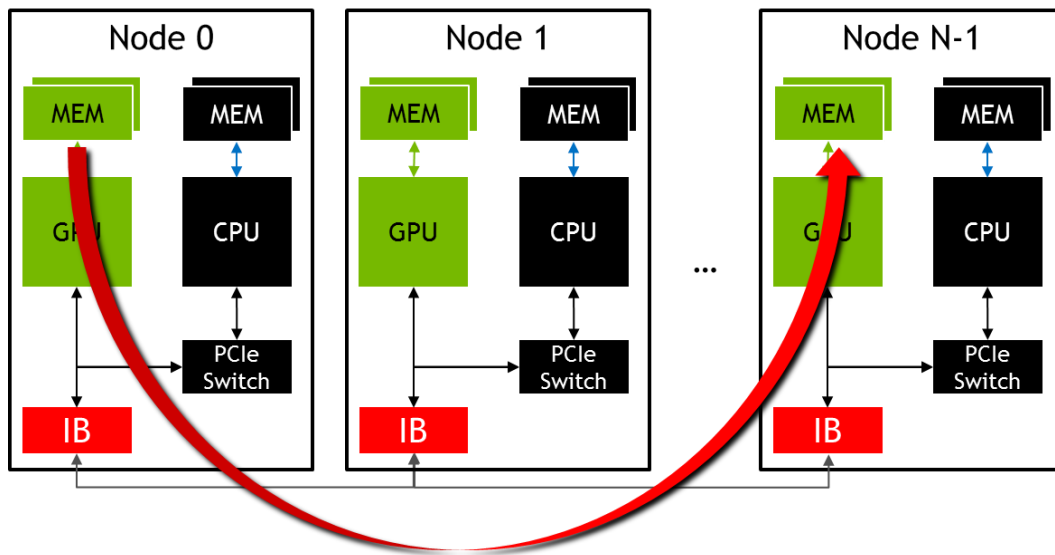
# MPI+CUDA

# MPI+CUDA

# MPI+CUDA



```
//MPI rank 0
MPI_Send(s_buf_d,size,MPI_CHAR,n-1,tag,MPI_COMM_WORLD);

//MPI rank n-1
MPI_Recv(r_buf_d,size,MPI_CHAR,0,tag,MPI_COMM_WORLD,&stat);
```

# YOU WILL LEARN

What MPI is

How to use MPI for inter GPU communication with CUDA and OpenACC

What CUDA-aware MPI is

What Multi Process Service is and how to use it

How to use NVIDIA tools in an MPI environment

How to hide MPI communication times

# MESSAGE PASSING INTERFACE - MPI

Standard to exchange data between processes via messages

Defines API to exchanges messages

Point to Point: e.g. `MPI_Send`, `MPI_Recv`

Collectives: e.g. `MPI_Reduce`

Multiple implementations (open source and commercial)

Bindings for C/C++, Fortran, Python, …

E.g. MPICH, OpenMPI, MVAPICH, IBM Platform MPI, Cray MPT, …
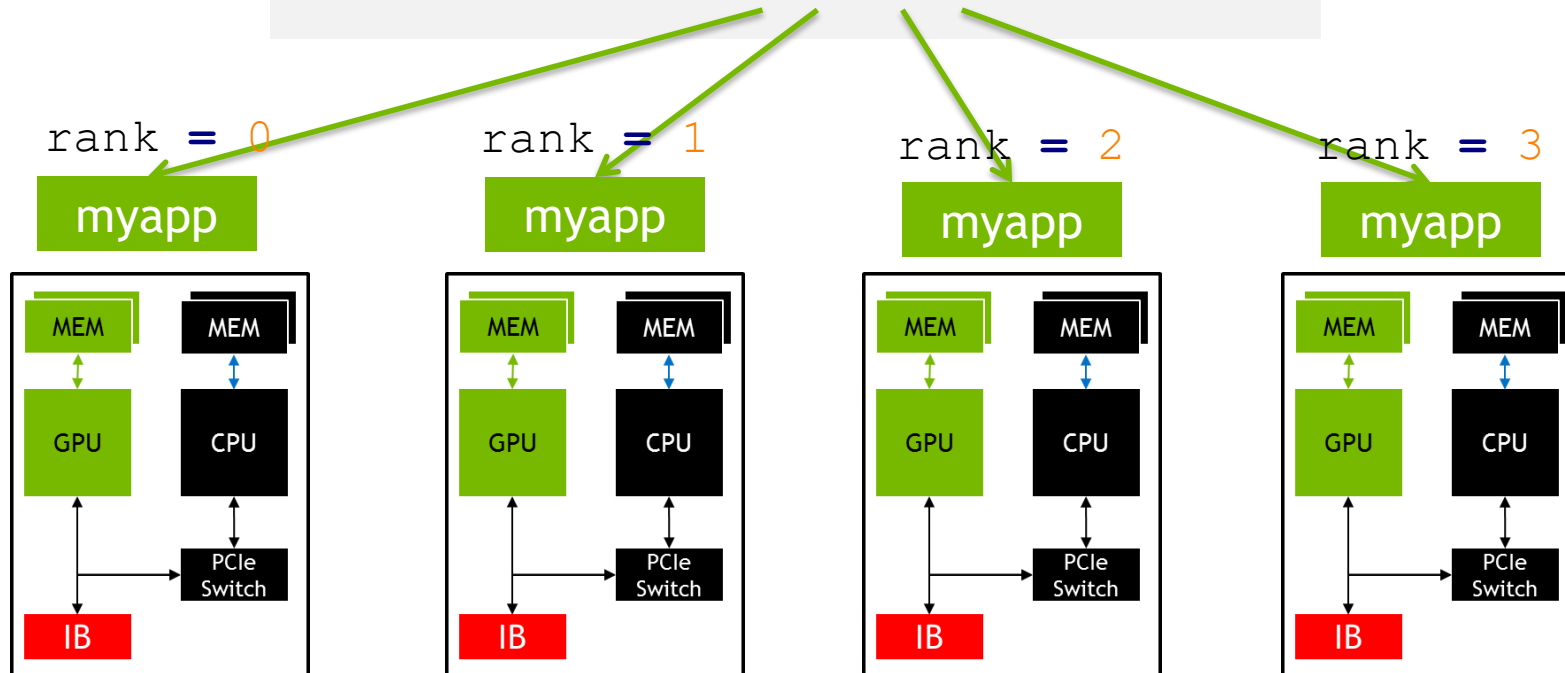
NVIDIA.

# MPI - SKELETON

```c
#include <mpi.h>
int main(int argc, char *argv[]) {
    int rank,size;
    /* Initialize the MPI library */
    MPI_Init(&argc,&argv);
    /* Determine the calling process rank and total number of ranks */
    MPI_Comm_rank(MPI_COMM_WORLD,&rank);
    MPI_Comm_size(MPI_COMM_WORLD,&size);
    /* Call MPI routines like MPI_Send, MPI_Recv, ... */
    ...
    /* Shutdown MPI library */
    MPI_Finalize();
    return 0;
}
```

# MPI

## Compiling and Launching

```
$ mpicc -o myapp myapp.c
$ mpirun -np 4 ./myapp <args>
```

rank = 0          rank = 1          rank = 2          rank = 3

myapp             myapp             myapp             myapp
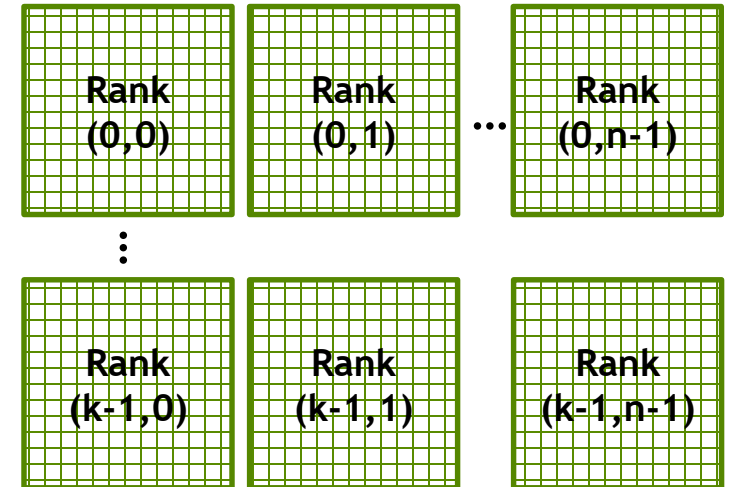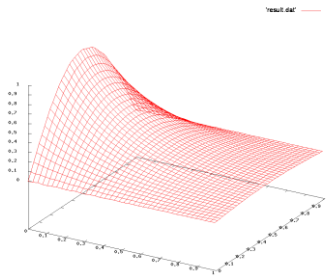
# A SIMPLE EXAMPLE

# EXAMPLE: JACOBI SOLVER

Solves the 2D-Laplace Equation on a rectangle

$$\Delta u(x, y) = 0 \ \forall \ (x, y) \in \Omega \backslash \delta\Omega$$

Dirichlet boundary conditions (constant values on boundaries)

$$u(x, y) = f(x, y) \in \delta\Omega$$

2D domain decomposition with n x k domains

| Rank (0,0) | Rank (0,1) | ... | Rank (0,n-1) |

⋮

| Rank (k-1,0) | Rank (k-1,1) | | Rank (k-1,n-1) |

NVIDIA.

# EXAMPLE: JACOBI SOLVER
## Single GPU

While not converged
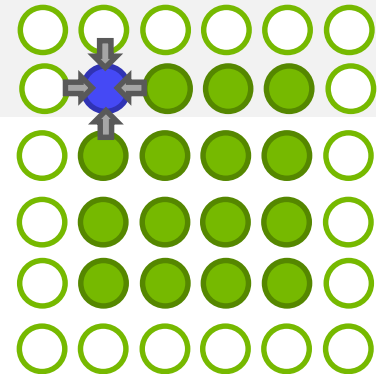
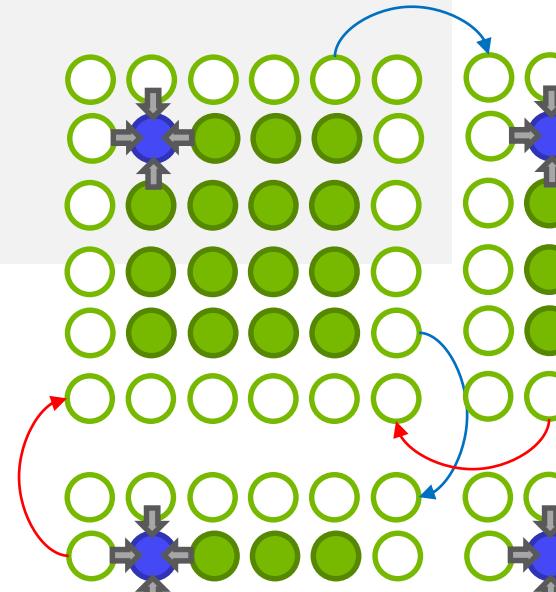Do Jacobi step:

```
for (int iy=1; iy < ny-1; ++iy)

for (int ix=1; ix < nx-1; ++ix)

    u_new[ix][iy] = 0.0f - 0.25f*( u[ix-1][iy] + u[ix+1][iy]

                                   + u[ix][iy-1] + u[ix][iy+1]);
```

Swap u_new and u

Next iteration

# EXAMPLE: JACOBI SOLVER
## Multi GPU

While not converged

Do Jacobi step:

```cpp
for (int iy=1; iy < ny-1; ++iy)

    for (int ix=1; ix < nx-1; ++ix)

        u_new[ix][iy] = 0.0f - 0.25f*( u[ix-1][iy] + u[ix+1][iy]

                                     + u[ix][iy-1] + u[ix][iy+1]);
```

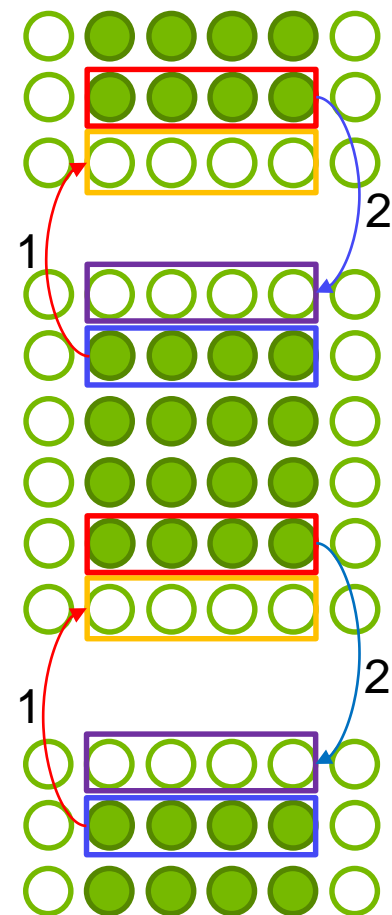Exchange halo with 2  4 neighbors

Swap `u_new` and `u`

Next iteration

# EXAMPLE JACOBI
## Top/Bottom Halo

```
MPI_Sendrecv(u_new+offset_first_row, m-2, MPI_DOUBLE, t_nb, 0,
             u_new+offset_bottom_boundary, m-2, MPI_DOUBLE, b_nb, 0,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);


MPI_Sendrecv(u_new+offset_last_row, m-2, MPI_DOUBLE, b_nb, 1,
             u_new+offset_top_boundary, m-2, MPI_DOUBLE, t_nb, 1,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
```
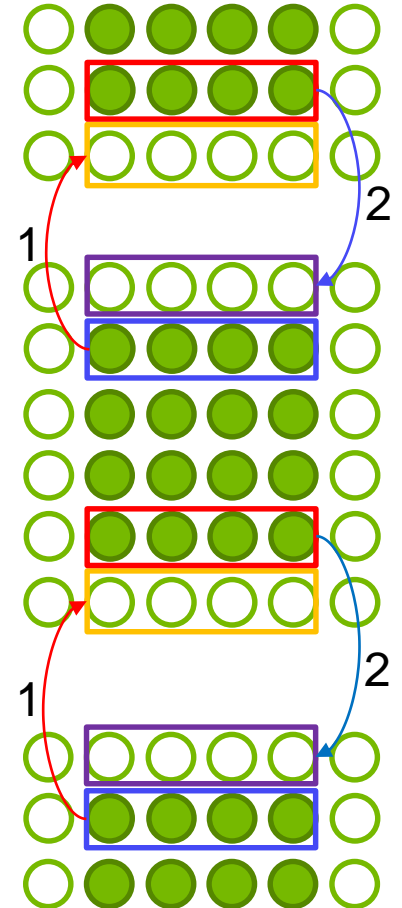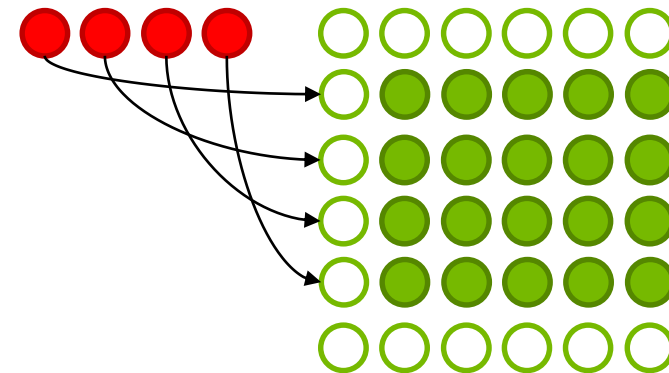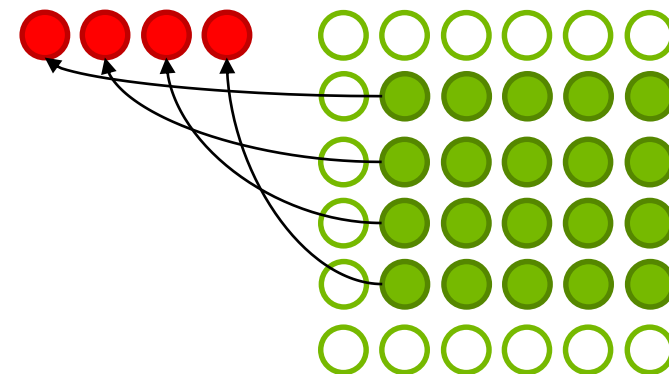
# EXAMPLE JACOBI

## Top/Bottom Halo



**OpenACC**

```
#pragma acc host_data use_device ( u_new ) {

MPI_Sendrecv(u_new+offset_first_row, m-2, MPI_DOUBLE, t_nb, 0,
             u_new+offset_bottom_boundary, m-2, MPI_DOUBLE, b_nb, 0,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
MPI_Sendrecv(u_new+offset_last_row, m-2, MPI_DOUBLE, b_nb, 1,
             u_new+offset_top_boundary, m-2, MPI_DOUBLE, t_nb, 1,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
}
```

**CUDA**

```
MPI_Sendrecv(u_new_d+offset_first_row, m-2, MPI_DOUBLE, t_nb, 0,
             u_new_d+offset_bottom_boundary, m-2, MPI_DOUBLE, b_nb, 0,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
MPI_Sendrecv(u_new_d+offset_last_row, m-2, MPI_DOUBLE, b_nb, 1,
             u_new_d+offset_top_boundary, m-2, MPI_DOUBLE, t_nb, 1,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
```

NVIDIA.

# EXAMPLE: JACOBI

## Left/Right Halo

```
//right neighbor omitted
#pragma acc parallel loop present ( u_new, to_left )
for ( int i=0; i<n-2; ++i )
        to_left[i] = u_new[(i+1)*m+1];


#pragma acc host_data use_device ( from_right, to_left ) {
  MPI_Sendrecv( to_left, n-2, MPI_DOUBLE, l_nb, 0,
                from_right, n-2, MPI_DOUBLE, r_nb, 0,
                MPI_COMM_WORLD, MPI_STATUS_IGNORE );
}


#pragma acc parallel loop present ( u_new, from_right )
for ( int i=0; i<n-2; ++i )
    u_new[(m-1)+(i+1)*m] = from_right[i];
```

OpenACC

# EXAMPLE: JACOBI

## Left/Right Halo

CUDA

```
//right neighbor omitted
pack<<<gs,bs,0,s>>>(to_left_d, u_new_d, n, m);
cudaStreamSynchronize(s);



  MPI_Sendrecv( to_left_d, n-2, MPI_DOUBLE, l_nb, 0,
                from_right_d, n-2, MPI_DOUBLE, r_nb, 0,
                MPI_COMM_WORLD, MPI_STATUS_IGNORE );



unpack<<<gs,bs,0,s>>>(u_new_d, from_right_d, n, m);
```

16

# LAUNCH MPI+CUDA/OPENACC PROGRAMS

Launch one process per GPU
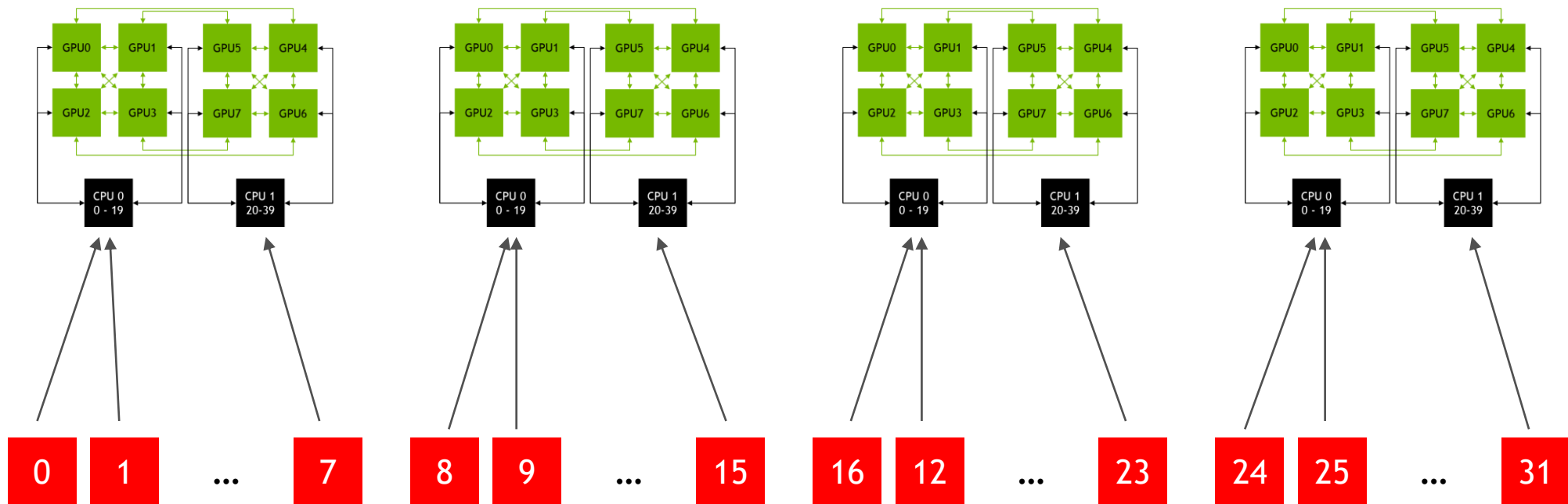
MVAPICH: `MV2_USE_CUDA`

```
$ MV2_USE_CUDA=1 mpirun -np ${np} ./myapp <args>
```

Open MPI: CUDA-aware features are enabled per default

Cray: `MPICH_RDMA_ENABLED_CUDA`

IBM Spectrum MPI:

```
$ mpirun -gpu -np ${np} ./myapp <args>
```

# HANDLING MULTIPLE MULTI GPU NODES

# HANDLING MULTIPLE MULTI GPU NODES

## How to determine the local rank? – MPI-3

```
MPI_Comm loc_comm;

MPI_Comm_split_type(MPI_COMM_WORLD, MPI_COMM_TYPE_SHARED, rank, MPI_INFO_NULL, &loc_comm);

int local_rank = -1;

MPI_Comm_rank(loc_comm,&local_rank);

MPI_Comm_free(&loc_comm);
```

# HANDLING MULTIPLE MULTI GPU NODES

# HANDLING MULTIPLE MULTI GPU NODES
## GPU-affinity

Use local rank:

```
int local_rank = -1;

MPI_Comm_rank(local_comm,&local_rank);

int num_devices = 0;

cudaGetDeviceCount(&num_devices);

cudaSetDevice(local_rank % num_devices);
```

NVIDIA.

# EXAMPLE JACOBI
## Top/Bottom Halo

without CUDA-aware MPI

**OpenACC**

```
#pragma acc update host(u_new[offset_first_row:m-2],u_new[offset_last_row:m-2])
MPI_Sendrecv(u_new+offset_first_row, m-2, MPI_DOUBLE, t_nb, 0,
             u_new+offset_bottom_boundary, m-2, MPI_DOUBLE, b_nb, 0,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
MPI_Sendrecv(u_new+offset_last_row, m-2, MPI_DOUBLE, b_nb, 1,
             u_new+offset_top_boundary, m-2, MPI_DOUBLE, t_nb, 1,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
#pragma acc update device(u_new[offset_top_boundary:m-2],u_new[offset_bottom_boundary:m-2])
```

**CUDA**

```
//send to bottom and receive from top top bottom omitted

cudaMemcpy(  u_new+offset_first_row,
             u_new_d+offset_first_row, (m-2)*sizeof(double), cudaMemcpyDeviceToHost);
MPI_Sendrecv(u_new+offset_first_row, m-2, MPI_DOUBLE, t_nb, 0,
             u_new+offset_bottom_boundary, m-2, MPI_DOUBLE, b_nb, 0,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
cudaMemcpy(  u_new_d+offset_bottom_boundary,
             u_new+offset_bottom_boundary, (m-2)*sizeof(double), cudaMemcpyDeviceToHost);
```

# THE DETAILS

# UNIFIED VIRTUAL ADDRESSING

*No UVA: Separate Address Spaces*      *UVA: Single Address Space*

# UNIFIED VIRTUAL ADDRESSING

**No UVA: Separate Address Spaces**

**UVA: Single Address Space**

One address space for all CPU and GPU memory

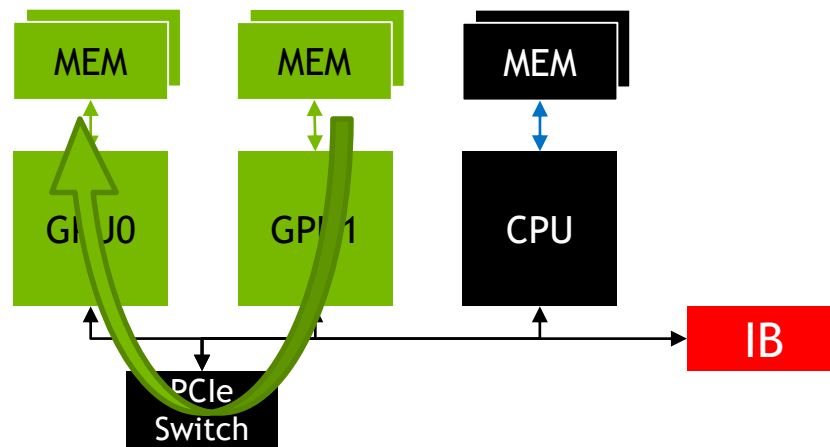   Determine physical memory location from a pointer value

   Enable libraries to simplify their interfaces (e.g. MPI and cudaMemcpy)

Supported on devices with compute capability 2.0+ for
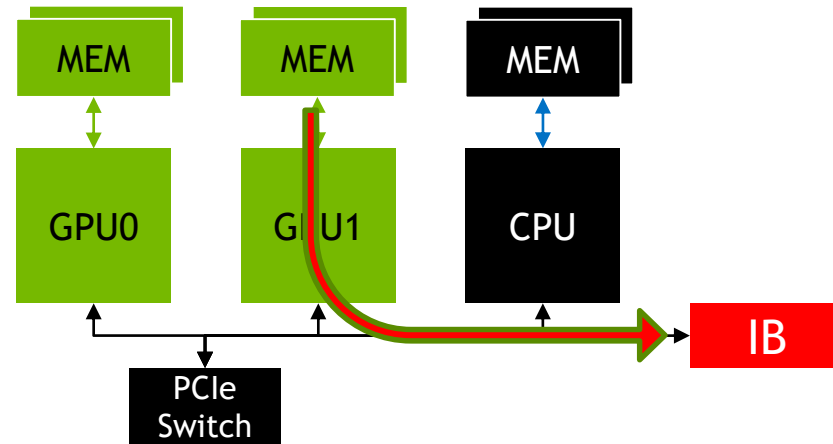
   64-bit applications on Linux and Windows (+TCC)

# NVIDIA GPUDIRECT™

## Peer to Peer Transfers

# NVIDIA GPUDIRECT™
## Support for RDMA

# CUDA-AWARE MPI

Example:

MPI Rank 0 `MPI_Send` from GPU Buffer

MPI Rank 1 `MPI_Recv` to GPU Buffer


Show how CUDA+MPI works in principle

Depending on the MPI implementation, message size, system setup, … situation might be different

Two GPUs in two nodes

# 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);
```

NVIDIA.

# MPI GPU TO REMOTE GPU

## Support for RDMA

# 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_Sendrecv

Time

# REGULAR MPI GPU TO REMOTE GPU



```
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);
```

NVIDIA.

# REGULAR MPI GPU TO REMOTE GPU

# PERFORMANCE RESULTS GPUDIRECT RDMA

## MVAPICH2-GDR 2.3a DGX-1V Tesla V100



Latency (1 Byte)    16.75 us    18.68 us    3.25 us

# PERFORMANCE RESULTS GPUDIRECT P2P

## MVAPICH2-GDR 2.3a DGX-1V Tesla V100

# MULTI PROCESS SERVICE (MPS) FOR MPI APPLICATIONS

# GPU ACCELERATION OF LEGACY MPI APPS

Typical legacy application

>   MPI parallel

>   Single or few threads per MPI rank (e.g. OpenMP)

Running with multiple MPI ranks per node

GPU acceleration in phases

>   Proof of concept prototype, …

>   Great speedup at kernel level

Application performance misses expectations

NVIDIA.

# MULTI PROCESS SERVICE (MPS)

## For Legacy MPI Applications



GPU parallelizable part
CPU parallel part
Serial part

With Hyper-Q/MPS
Available on Tesla/Quadro with CC 3.5+
(e.g. K20, K40, K80, M40,...)

N=1     N=2     N=4     N=8          N=1     N=2     N=4     N=8

Multicore CPU only                    GPU-accelerated

# PROCESSES SHARING GPU WITHOUT MPS

## No Overlap

# PROCESSES SHARING GPU WITHOUT MPS

## Context Switch Overhead



Context Switch

# PROCESSES SHARING GPU WITH MPS

## Maximum Overlap



Process A
Context A

Process B
Context B

MPS Process

GPU

Kernels from Process A

Kernels from Process B

# PROCESSES SHARING GPU WITH MPS

## No Context Switch Overhead



NVIDIA.

# HYPER-Q/MPS CASE STUDY: UMT



Enables overlap between copy and compute of different processes

GPU sharing between MPI ranks increases utilization

# HYPER-Q/MPS CASE STUDIES
## CPU Scaling Speedup

# HYPER-Q/MPS CASE STUDIES

## Additional Speedup with MPS



46 NVIDIA.

# USING MPS

No application modifications necessary

Not limited to MPI applications

MPS control daemon

    Spawn MPS server upon CUDA
    application startup

```
#Typical Setup

nvidia-smi -c EXCLUSIVE_PROCESS

nvidia-cuda-mps-control -d


#On Cray XK/XC systems

export CRAY_CUDA_MPS=1
```

# MPS: IMPROVEMENTS WITH VOLTA

**More MPS clients per GPU:** 48 instead of 16

**Less overhead:** Volta MPS clients submit work directly to the GPU without passing through the MPS server.

**More security:** Each Volta MPS client owns its own GPU address space instead of sharing GPU address space with all other MPS clients.

**More control:** Volta MPS supports limited execution resource provisioning for Quality of Service (QoS). ->
`CUDA_MPS_ACTIVE_THREAD_PERCENTAGE`

# MPS SUMMARY

Easy path to get GPU acceleration for legacy applications

Enables overlapping of memory copies and compute between different MPI ranks

Remark: MPS adds some overhead!

# DEBUGGING AND PROFILING

# TOOLS FOR MPI+CUDA APPLICATIONS

Memory checking: `cuda-memcheck`

Debugging: `cuda-gdb`

Profiling: `nvprof` and the NVIDIA Visual Profiler (`nvvp`)

# MEMORY CHECKING WITH CUDA-MEMCHECK

`cuda-memcheck` is a tool similar to Valgrind's memcheck

Can be used in a MPI environment

```
mpiexec -np 2 cuda-memcheck ./myapp <args>
```

Problem: Output of different processes is interleaved

Solution: Use save or log-file command line options

OpenMPI: `OMPI_COMM_WORLD_RANK`

MVAPICH2: `MV2_COMM_WORLD_RANK`

```
mpirun -np 2 cuda-memcheck                              \

  --log-file name.%q{OMPI_COMM_WORLD_RANK}.log          \

  --save name.%q{OMPI_COMM_WORLD_RANK}.memcheck         \

  ./myapp <args>
```

NVIDIA.

# MEMORY CHECKING WITH CUDA-MEMCHECK

# MEMORY CHECKING WITH CUDA-MEMCHECK

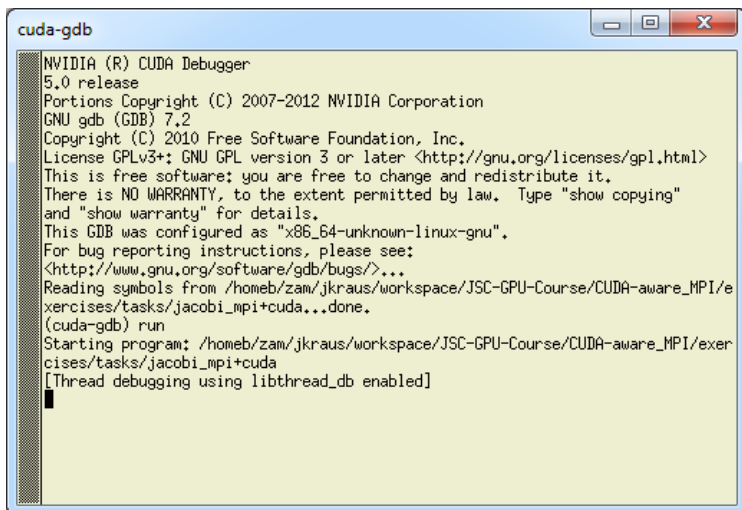## Read Output Files with `cuda-memcheck --read`

# DEBUGGING MPI+CUDA APPLICATIONS
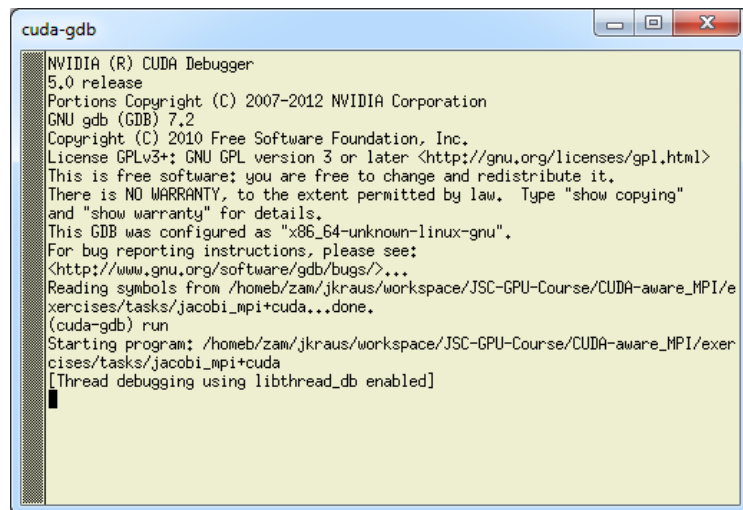
## Using `cuda-gdb` with MPI Applications

Use `cuda-gdb` **just like** `gdb`

For smaller applications, just launch xterms and `cuda-gdb`

```
mpiexec -x -np 2 xterm -e cuda-gdb ./myapp <args>
```



⬢ **NVIDIA.**

# DEBUGGING MPI+CUDA APPLICATIONS
## cuda-gdb Attach

```
if ( rank == 0 ) {
        int i=0;
        printf("rank %d: pid %d on %s ready for attach\n.", rank, getpid(),name);
        while (0 == i) { sleep(5); }
}


> mpiexec -np 2 ./jacobi_mpi+cuda
Jacobi relaxation Calculation: 4096 x 4096 mesh with 2 processes and one Tesla M2070 for
each process (2049 rows per process).
rank 0: pid 30034 on judge107 ready for attach
> ssh judge107
jkraus@judge107:~> cuda-gdb --pid 30034
```

NVIDIA.

# DEBUGGING MPI+CUDA APPLICATIONS

## CUDA_DEVICE_WAITS_ON_EXCEPTION



⊗ nVIDIA.

# DEBUGGING MPI+CUDA APPLICATIONS

With `CUDA_ENABLE_COREDUMP_ON_EXCEPTION=1` core dumps are generated in case of an exception:

    Can be used for offline debugging

    Helpful if live debugging is not possible

`CUDA_ENABLE_CPU_COREDUMP_ON_EXCEPTION:` Enable/Disable CPU part of core dump (enabled by default)

`CUDA_COREDUMP_FILE:` Specify name of core dump file

Open GPU:     `(cuda-gdb) target cudacore core.cuda`

Open CPU+GPU: `(cuda-gdb) target core core.cpu core.cuda`

# DEBUGGING MPI+CUDA APPLICATIONS

## CUDA_ENABLE_COREDUMP_ON_EXCEPTION

# DEBUGGING MPI+CUDA APPLICATIONS

## CUDA_ENABLE_COREDUMP_ON_EXCEPTION



NVIDIA.

# DEBUGGING MPI+CUDA APPLICATIONS
## Third Party Tools

Allinea DDT debugger

Rogue Wave TotalView

# PROFILING MPI+CUDA APPLICATIONS

Using `nvprof`+NVVP

New since CUDA 9

Embed MPI rank in output filename, process name, and context name (OpenMPI)

```
mpirun -np $np nvprof --output-profile profile.%q{OMPI_COMM_WORLD_RANK} \
                      --process-name "rank %q{OMPI_COMM_WORLD_RANK}"     \
                      --context-name "rank %q{OMPI_COMM_WORLD_RANK}"     \
                      --annotate-mpi openmpi
```

MVAPICH2: `MV2_COMM_WORLD_RANK`
                `--annotate-mpi mpich`

Alternatives:

Only save the textual output (`--log-file`)

Collect data from all processes that run on a node (`--profile-all-processes`)

NVIDIA.

# PROFILING MPI+CUDA APPLICATIONS
## Using `nvprof`+NVVP



**NVIDIA.**

# PROFILING MPI+CUDA APPLICATIONS
## Using `nvprof`+NVVP



NVIDIA.

# PROFILING MPI+CUDA APPLICATIONS
## Third Party Tools

Multiple parallel profiling tools are CUDA-aware

Score-P

Vampir

Tau

These tools are good for discovering MPI issues as well as basic CUDA performance inhibitors.

# ADVANCED MPI ON GPUS

# BEST PRACTICE: USE NON-BLOCKING MPI

**BLOCKING**

```
#pragma acc host_data use_device ( u_new ) {
MPI_Sendrecv(u_new+offset_first_row, m-2, MPI_DOUBLE, t_nb, 0,
             u_new+offset_bottom_boundary, m-2, MPI_DOUBLE, b_nb, 0,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
MPI_Sendrecv(u_new+offset_last_row, m-2, MPI_DOUBLE, b_nb, 1,
             u_new+offset_top_boundary, m-2, MPI_DOUBLE, t_nb, 1,
             MPI_COMM_WORLD, MPI_STATUS_IGNORE);
}
```

**NON-BLOCKING**

```
MPI_Request t_b_req[4];
#pragma acc host_data use_device ( u_new ) {
  MPI_Irecv(u_new+offset_top_boundary,m-2,MPI_DOUBL
  MPI_Irecv(u_new+offset_bottom_boundary,m-2,MPI          _req+1);
  MPI_Isend(u_new+offset_last_row,m-2,MPI_DOUBLE           _b_req+2);
  MPI_Isend(u_new+offset_first_row,m-2,MPI_DOUBLE          MM_WORLD,t_b_req+3);
}
MPI_Waitall(4, t_b_req, MPI_STATUSES_IGNORE);
```

Gives MPI more opportunities to build efficient piplines

NVIDIA.

# COMMUNICATION + COMPUTATION OVERLAP

## OpenMPI 3.0.1RC1 – DGX-1V – 2048x2048

# COMMUNICATION + COMPUTATION OVERLAP

No Overlap

Process Whole Domain | MPI

Overlap

Boundary and inner domain processing can overlap

Process inner domain | Possible gain

Process boundary domain | Dependency | MPI

NVIDIA.

# COMMUNICATION + COMPUTATION OVERLAP
## CUDA with Streams

```
process_boundary_and_pack<<<gs_b,bs_b,0,s1>>>(u_new_d,u_d,to_left_d,to_right_d,n,m);


process_inner_domain<<<gs_id,bs_id,0,s2>>>(u_new_d, u_d,to_left_d,to_right_d,n,m);


cudaStreamSynchronize(s1);        //wait for boundary
MPI_Request req[8];

//Exchange halo with left, right, top and bottom neighbor

MPI_Waitall(8, req, MPI_STATUSES_IGNORE);
unpack<<<gs_s,bs_s,0,s2>>>(u_new_d, from_left_d, from_right_d, n, m);


cudaDeviceSynchronize();        //wait for iteration to finish
```

# COMMUNICATION + COMPUTATION OVERLAP
## OpenACC with Async Queues

```
#pragma acc parallel loop present ( u_new, u, to_left, to_right ) async(1)
for ( ... )
    //Process boundary and pack to_left and to_right
#pragma acc parallel loop present ( u_new, u ) async(2)
for ( ... )
    //Process inner domain
#pragma acc wait(1)                      //wait for boundary
MPI_Request req[8];
#pragma acc host_data use_device ( from_left, to_left, form_right, to_right, u_new ) {
    //Exchange halo with left, right, top and bottom neighbor
}
MPI_Waitall(8, req, MPI_STATUSES_IGNORE);
#pragma acc parallel loop present ( u_new, from_left, from_right ) async(2)
for ( ... )
    //unpack from_left and from_right
#pragma acc wait                         //wait for iteration to finish
```

NVIDIA.

# COMMUNICATION + COMPUTATION OVERLAP

## OpenMPI 3.0.1RC1 – DGX-1V – 2048x2048



NVIDIA.

# HIGH PRIORITY STREAMS

Improve scalability with high priority streams

__host__ cudaError_t cudaStreamCreateWithPriority ( cudaStream_t* pStream, unsigned int flags, int priority )

Use-case: MD Simulations

| Stream 1 | Comp. Local Forces |
| Stream 2 | Ex. Non-local Atom pos. → Comp. Non-Local Forces → Ex. Non-local forces |
| Stream 1 (LP) | Comp. Local Forces → Possible gain |
| Stream 2 (HP) | Ex. Non-local Atom pos. → Comp. Non-Local Forces → Ex. Non-local forces |

# MPI AND UNIFIED MEMORY
## CAVEAT

Using Unified Memory with a non Unified Memory-aware MPI might fail with errors or even worse silently produce wrong results, e.g. when registering Unified Memory for RDMA.

**Use a Unified Memory-aware MPI,**

**e.g. OpenMPI since 1.8.5 or MVAPICH2-GDR since 2.2b**

Unified Memory-aware: CUDA-aware MPI with support for Unified Memory

NVIDIA.

# MPI AND UNIFIED MEMORY

## Performance Implications

Unified Memory can be used by any processor in the system

Memory pages of a Unified Memory allocation may migrate between processors memories to ensure coherence and maximize performance

Different data paths are optimal for performance depending on where the data is: e.g. NVLink between peer GPUs

The MPI implementation needs to know where the data is,

but it can't!

# MPI AND UNIFIED MEMORY
## Performance Implications - Simple Example

```
cudaMallocManaged( &array, n*sizeof(double), cudaMemAttachGlobal );


while( ... ) {

        foo(array,n);

        MPI_Send(array,...);

        foo(array,n);

}
```

**NVIDIA.**

# MPI AND UNIFIED MEMORY

## Performance Implications - Simple Example

▸ If foo is a CPU function pages of array might migrate to System Memory

▸ If foo is a GPU function pages of array might migrate to GPU Memory

▸ The MPI implementation is not aware of the application and thus doesn't know where array is and what's optimal

```
while( ... ) {

        foo(array,n);

        MPI_Send(array,...);

        foo(array,n);

}
```

NVIDIA.

# MPI AND UNIFIED MEMORY
## The Future with Data Usage Hints

Tell where the application intends to use the data

```
cudaMallocManaged( &array, n*sizeof(double), cudaMemAttachGlobal );

cudaMemAdvise(array,n*sizeof(double),cudaMemAdviseSetPreferredLocation,device);

while( ... ) {

        foo(array,n);

        MPI_Send(array,...);

        foo(array,n);

}
```

Array is intended to be used on the GPU with the id device

Remark: Data Usage Hints are available since CUDA 8, but currently not evaluated by any Unified Memory-aware MPI implementation.

NVIDIA.

# MPI AND UNIFIED MEMORY
## The Future with Data Usage Hints

Tell where the application intends to use the data

```
cudaMallocManaged( &array, n*sizeof(double), cudaMemAttachGlobal );

cudaMemAdvise(array,n*sizeof(double),cudaMemAdviseSetPreferredLocation, cudaCpuDeviceId);

while( ... ) {

        foo(array,n);

        MPI_Send(array,...);

        foo(array,n);

}
```

**Array is intended to be used on the CPU**

Remark: Data Usage Hints are available since CUDA 8, but currently not evaluated by any Unified Memory-aware MPI implementation.

⊚ **NVIDIA.**

# MPI AND UNIFIED MEMORY
## The Future with Data Usage Hints - Summary

Data usage hints can be queried by the MPI Implementation and allow it to take the optimal data path

If the application lies about the data usage hints it will run correctly but performance will be affected

Performance tools help to identify missing or wrong data usage hints

Data usage hints are general useful for the Unified Memory system and can improve application performance.

Remark: Data Usage Hints are only hints to guide the data usage policies of the Unified Memory system. The Unified Memory system might ignore them, e.g. to ensure coherence or in oversubscription scenarios.

# MPI AND UNIFIED MEMORY
## Current Status

Available Unified Memory-aware MPI implementations

- OpenMPI (since 1.8.5)

- MVAPICH2-GDR (since 2.2b)

  - Performance improvements with 2.2RC1 for Intranode GPU to GPU communication

Currently both don't evaluate Data Usage Hints, i.e. all Unified Memory is treated as Device Memory

Good performance if all buffers used in MPI are touched mainly on the GPU.

# MPI AND UNIFIED MEMORY
## Without Unified Memory-aware MPI

Only use non Unified Memory Buffers for MPI: cudaMalloc, cudaMallocHost or malloc

Application managed non Unified Memory Buffers also allow to work around current missing cases in Unified Memory-aware MPI Implementations.

**NVIDIA.**

# DETECTING CUDA-AWARENESS

OpenMPI (since 2.0.0):

Macro:

`MPIX_CUDA_AWARE_SUPPORT`

Function for runtime decisions

`MPIX_Query_cuda_support()`

Include `mpi-ext.h` for both.

See http://www.open-mpi.org/faq/?category=runcuda#mpi-cuda-aware-support

NVIDIA.