# DEVICE-INITIATED COMMUNICATION WITH NVSHMEM

JIRI KRAUS, PRINCIPAL DEVTECH COMPUTE

# CPU-INITIATED COMMUNICATION
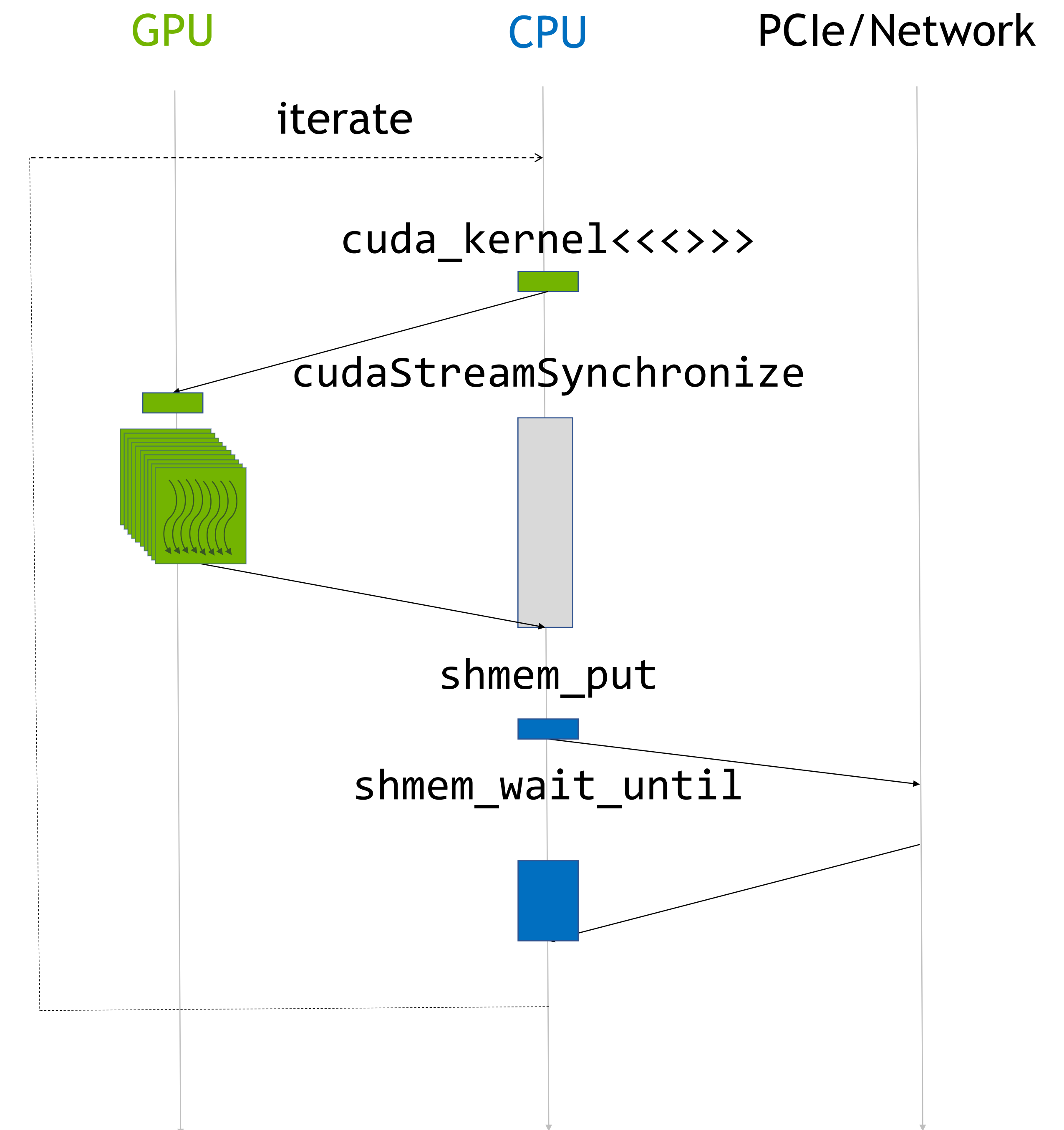
- Compute on GPU

- Communication from CPU

Synchronization at boundaries

Commonly used model, but –

- Offload latencies in critical path

- Communication is not overlapped

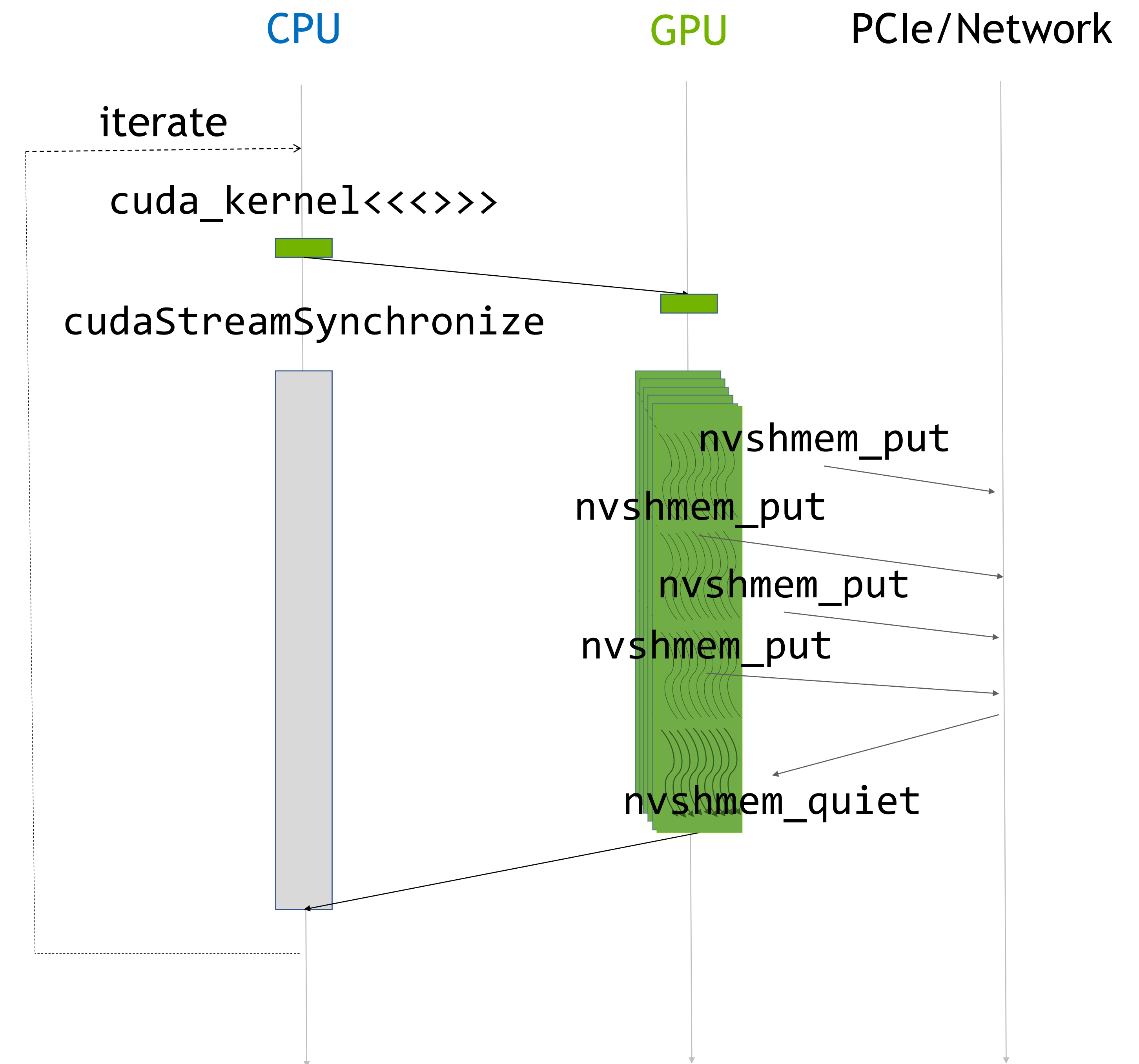Hiding increased code complexity, not hiding limits strong scaling

GPU          CPU          PCIe/Network

iterate

cuda_kernel<<<>>>

cudaStreamSynchronize

shmem_put

shmem_wait_until

# GPU-INITIATED COMMUNICATION
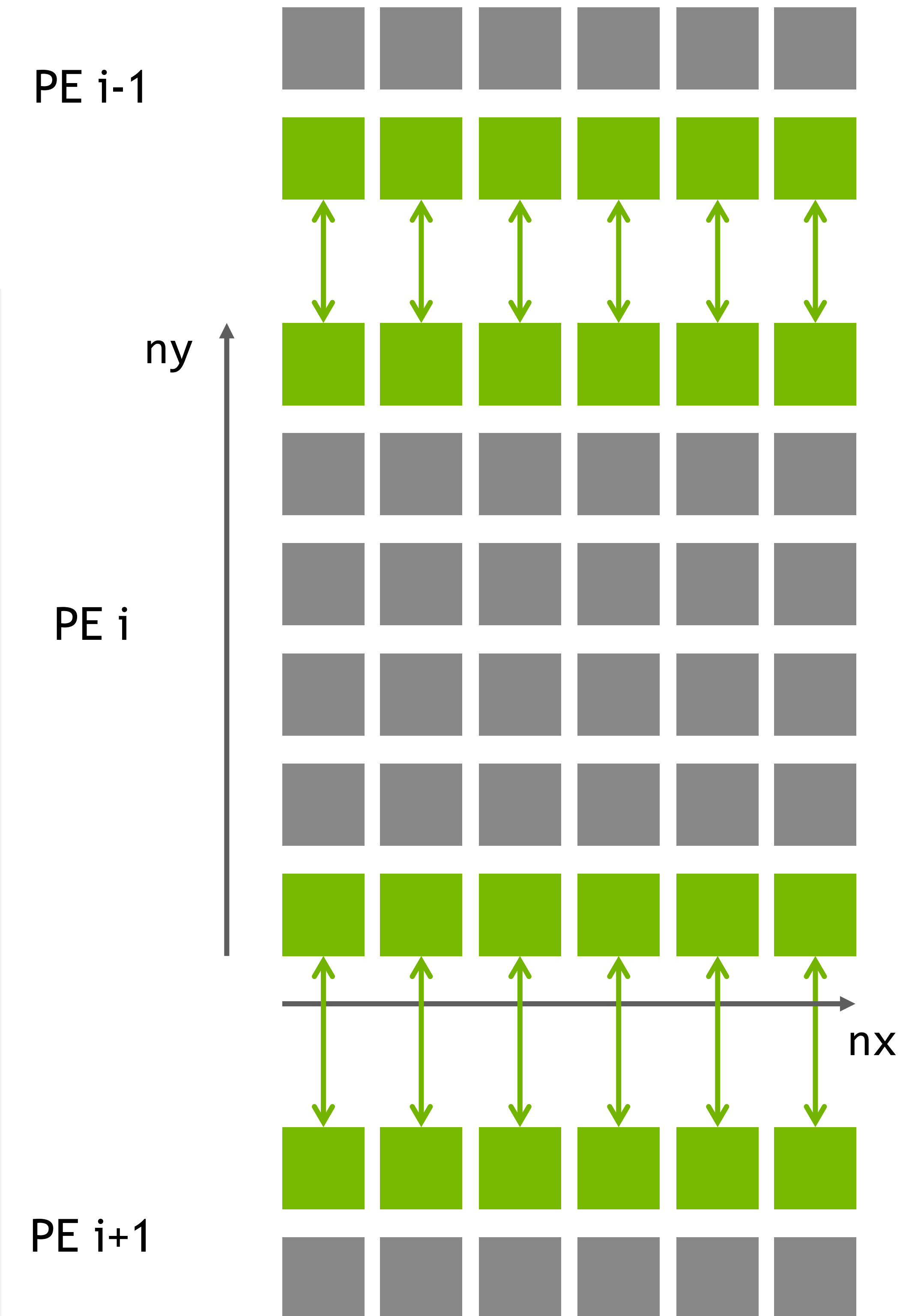
- Compute on GPU

- Communication from GPU

Benefits

- Eliminates offloads latencies

- Compute and communication overlap by threading

- Easier to express algorithms with inline communication

CPU     GPU     PCIe/Network

iterate

cuda_kernel<<<>>>

cudaStreamSynchronize

nvshmem_put

nvshmem_put

nvshmem_put

nvshmem_put

nvshmem_quiet

Slide 3

# THREAD-LEVEL COMMUNICATION

- Allows fine grained communication and computation overlap

- Efficient mapping to NVLink fabric on DGX systems

```
__global__ void stencil_single_step(float *u, float *v, …) {
  int ix = get_ix(blockIdx, blockDim, threadIdx);
  int iy = get_iy(blockIdx, blockDim, threadIdx);
  compute(u, v, ix, iy);
  // Thread-level data communication API

  if (iy == 1)
    nvshmem_float_p(u+(ny+1)*nx+ix, u[nx+ix], top_pe);
  if (iy == ny)
    nvshmem_float_p(u+ix, u[ny*nx+ix], bottom_pe);
}
for (int iter = 0; iter < N; iter++) {
  swap(u, v);
  stencil_single_step<<<..., stream>>>(u, v, …);
  nvshmem_barrier_all_on_stream(stream);
}
```
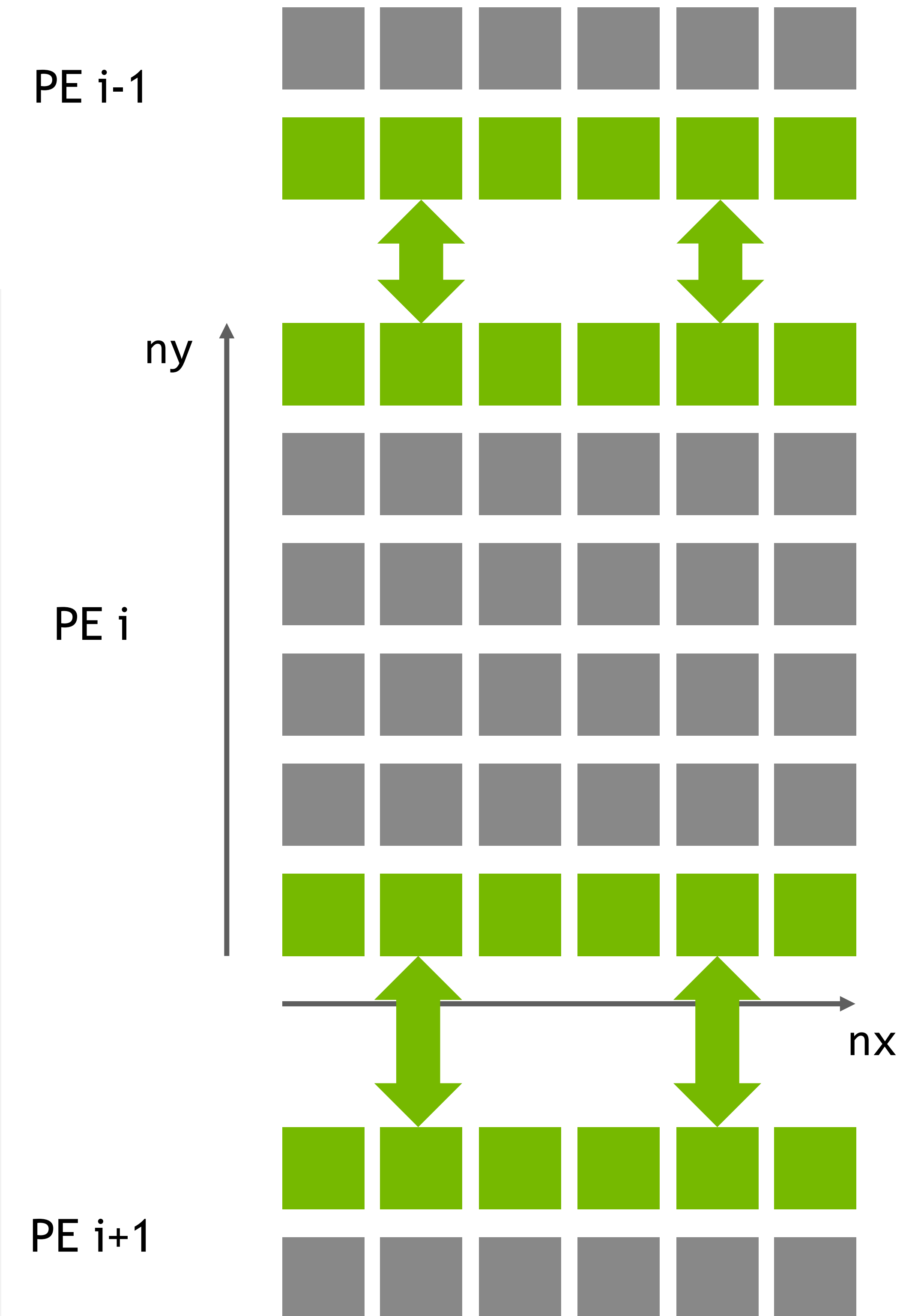
PE i-1

ny

PE i

PE i+1

nx

NVIDIA

# THREAD-GPOUP COMMUNICATION

- NVSHMEM operations can be issued by all threads in a block/warp

- More efficient data transfers over networks like IB

- Still allows inter-warp/inter-block overlap

```
__global__ void stencil_single_step(float *u, float *v, …) {
  int ix = get_ix(blockIdx, blockDim, threadIdx);
  int iy = get_iy(blockIdx, blockDim, threadIdx);
  compute(u, v, ix, iy);
  // Thread block-level communication API
  int boffset = get_block_offet(blockIdx,blockDim);
  if (blockIdx.y == 0)
    nvshmemx_float_put_nbi_block(u+(ny+1)*nx+boffset, u+nx+boffset, blockDim.x, top_pe);
  if (blockIdx.y == (blockDim.y-1))
    nvshmemx_float_put_nbi_block(u+boffset, u+ny*nx+boffset, blockDim.x, bottom_pe);
}
for (int iter = 0; iter < N; iter++) {
  swap(u, v);
  stencil_single_step<<<..., stream>>>(u, v, …);
  nvshmem_barrier_all_on_stream(stream);
}
```
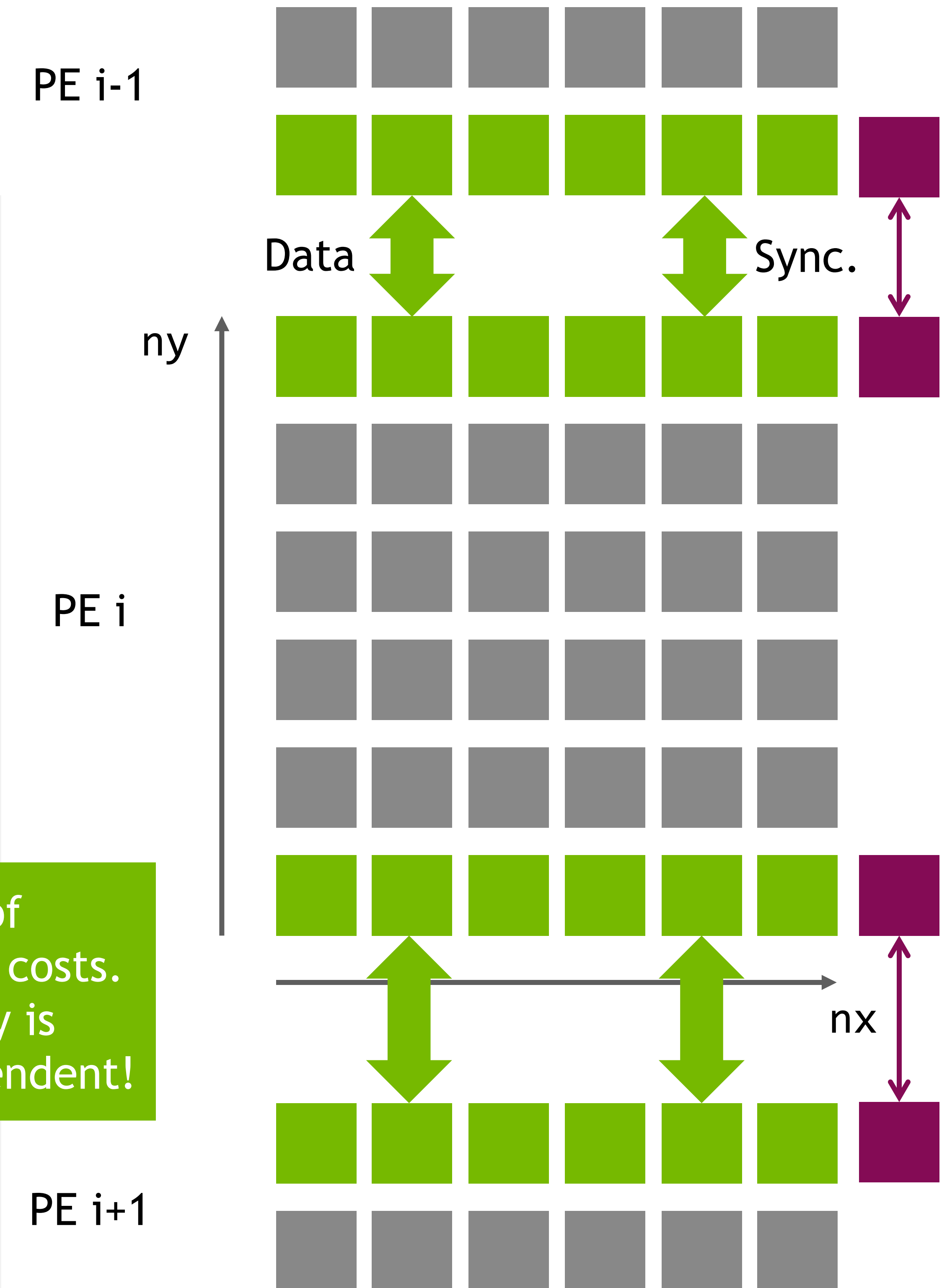
# IN-KERNEL SYNCHRONIZATION

- Point-to-point synchronization across PEs within a kernel

- Enables kernel fusion

```
__global__ void stencil_multi_step(float *u, float *v, int N, int *sync, …) {
  int ix = get_ix(blockIdx, blockDim, threadIdx);
  int iy = get_iy(blockIdx, blockDim, threadIdx);
  for (int iter = 0; iter < N; iter++) {
    swap(u, v); compute(u, v, ix, iy);
    // Thread block-level data exchange (assume even/odd iter buffering)
    int boffset = get_block_offet(blockIdx,blockDim);
    if (blockIdx.y == 0)
      nvshmemx_float_put_nbi_block(u+(ny+1)*nx+boffset, u+nx+boffset, blockDim.x, top_pe);
    if (blockIdx.y == (blockDim.y-1))
      nvshmemx_float_put_nbi_block(u + boffset, u+ny*nx+boffset,   blockDim.x, bottom_pe);
    if (blockIdx.y == 0 || blockIdx.y == (blockDim.y-1)) {
      __syncthreads();
      nvshmem_quiet();
      if (threadIdx.x == 0 && threadIdx.y == 0) {
        nvshmem_atomic_inc(sync, top_pe);
        nvshmem_atomic_inc(sync, bottom_pe);
    }}
    nvshmem_wait_until(sync, NVSHMEM_CMP_GT, 2*iter*gridDim.x);
}}
```

Be aware of synchronization costs. Best strategy is application dependent!



PE i-1

Data    Sync.

ny

PE i

nx

PE i+1

# COLLECTIVE KERNEL LAUNCH

Ensures progress when using device-side inter-kernel synchronization

| NVSHMEM Usage | CUDA Kernel launch |
|---|---|
| Device-Initiated Communication | Execution config syntax <<<...>>> or launch APIs |
| Device-Initiated Synchronization | `nvshmemx_collective_launch` |

- CUDA's throughput computing model allows (encourages) grids much larger than a GPU can fit

- Inter-kernel synchronization requires producer and consumer threads to execute concurrently

- Collective launch guarantees co-residency using CUDA cooperative launch and requirement of 1PE/GPU
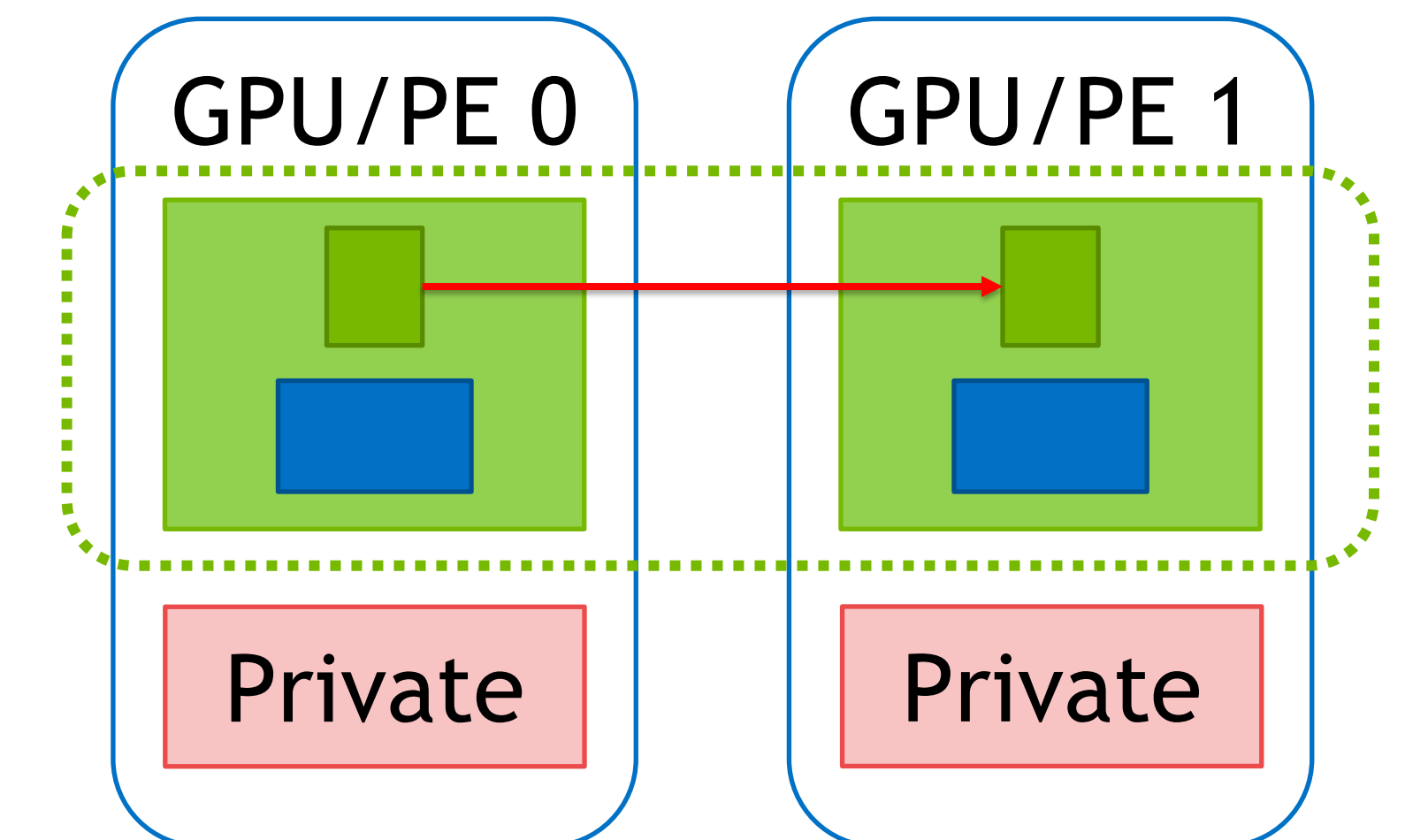
nVIDIA

# NVSHMEM API
## single element put

```
__device__ void nvshmem_TYPENAME_p(TYPE *dest, TYPE value, int pe)
```

- dest [OUT]: Symmetric address of the destination data object.

- value [IN]: The value to be transferred to dest.

- pe [IN]: The number of the remote PE.

See: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/rma.html#nvshmem-p

TYPENAME can be: float, double, char, schar, short, int, long, longlong, uchar, ushort, uint, ulong, ulonglong, …, ptrdiff
(see: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/rma.html#stdrmatypes)

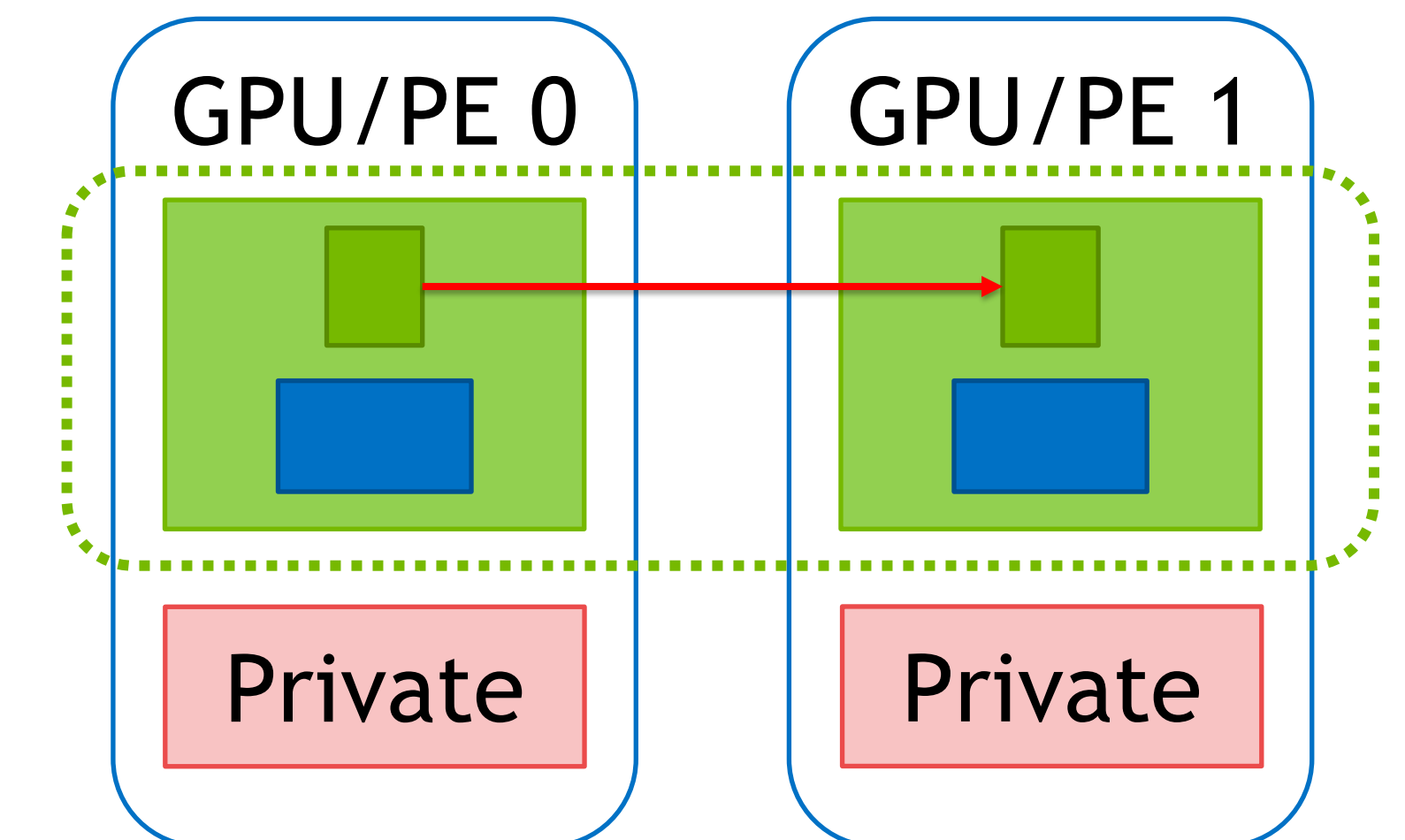GPU/PE 0    GPU/PE 1

Private    Private

# NVSHMEM API
## nonblocking block cooperative put

```
__device__ void nvshmemx_TYPENAME_put_nbi_block(TYPE *dest, const TYPE *source, size_t nelems, int pe)
```

- `dest` [OUT]: Symmetric address of the destination data object.

- `source` [IN]: Symmetric address of the object containing the data to be copied.

- `nelems` [IN]: Number of elements in the `dest` and `source` arrays.

- `pe` [IN]: The number of the remote PE.

Cooperative call: Needs to be called by all threads in a block. thread and warp are also available.

x in nvshmem**x** marks API as extension of the OpenSHMEM APIs.

See: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/rma.html?highlight=nvshmemx_typename_put_nbi_block#nvshmem-put-nbi

TYPENAME can be: float, double, char, schar, short, int, long, longlong, uchar, ushort, uint, ulong, ulonglong, …, ptrdiff
(see: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/rma.html#stdrmatypes)

GPU/PE 0    GPU/PE 1

Private    Private

# NVSHMEM API
## ordering and completion

```
__device__ void nvshmem_quiet(void)
```

Ensures completion of all operations on symmetric data objects issued by the calling PE.

See: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/ordering.html#nvshmem-quiet

# NVSHMEM API
signal operation

```
__device__ inline void nvshmemx_signal_op(uint64_t *sig_addr, uint64_t signal, int sig_op, int pe)
```

- sig_addr [OUT]: Symmetric address of the signal word to be updated.

- signal [IN]: The value used to update sig_addr.

- sig_op [IN]: Operation used to update sig_addr with signal. (NVSHMEM_SIGNAL_SET or NVSHMEM_SIGNAL_ADD)

- pe [IN]: The number of the remote PE.

x in nvshmemx marks API as extension of the OpenSHMEM APIs.

See: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/signal.html#nvshmemx-signal-op

NVIDIA

# NVSHMEM API
atomic operation

```
__device__ void nvshmem_TYPENAME_atomic_inc(TYPE *dest, int pe)
```

- dest [OUT]: Symmetric address of the signal word to be updated.

- pe [IN]: The number of the remote PE.

These routines perform an atomic increment operation on the dest data object on PE.

See: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/amo.html#nvshmem-atomic-inc

TYPENAME can be: float, double, char, schar, short, int, long, longlong, uchar, ushort, uint, ulong, ulonglong, …, ptrdiff
(see: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/rma.html#stdrmatypes)

# NVSHMEM API
wait operations

```
__device__ void nvshmem_TYPENAME_wait_until_all(TYPE *ivars, size_t nelems, const int *status,

                                                int cmp, TYPE cmp_value)



__device__ void nvshmem_TYPENAME_wait_until(TYPE *ivar, int cmp, TYPE cmp_value)
```

- ivars|ivar [IN]: Symmetric address of an array of remotely accessible data objects.|Symmetric address of a remotely accessible data object.

- nelems [IN]: The number of elements in the ivars array.

- status [IN]: Local address of an optional mask array of length nelems that indicates which elements in ivars are excluded from the wait set. Set to NULL when not used.

- cmp [IN]: A comparison operator (NVSHMEM_CMP_EQ, NVSHMEM_CMP_NE, NVSHMEM_CMP_GT, NVSHMEM_CMP_GE, NVSHMEM_CMP_LT, NVSHMEM_CMP_LE) that compares elements of ivars|ivar with cmp_value.

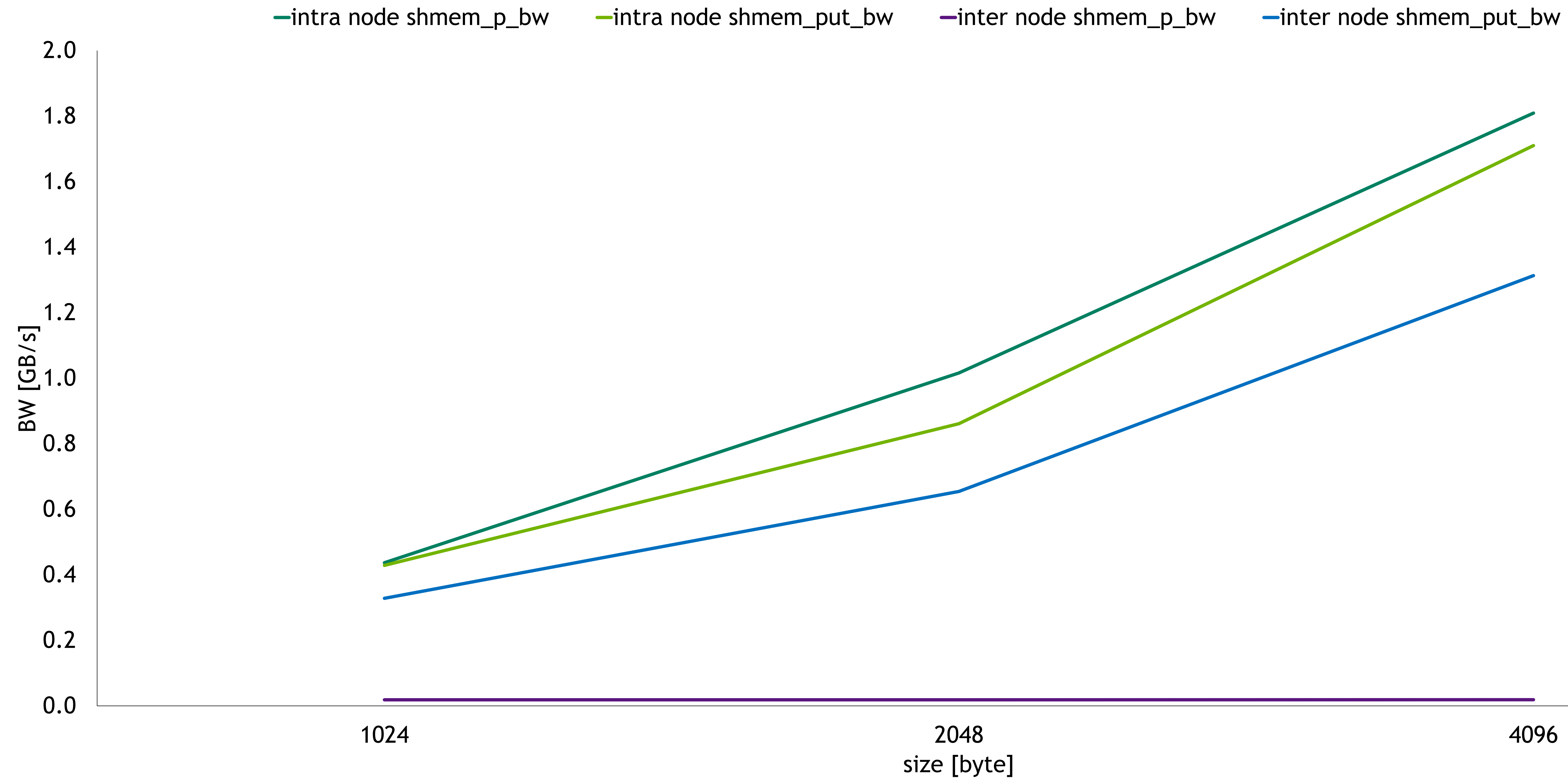- cmp_value [IN]: The value to be compared with the objects pointed to by ivars.

See: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/sync.html#nvshmem-wait-until-all and
https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/sync.html#nvshmem-wait-until

TYPENAME can be: float, double, char, schar, short, int, long, longlong, uchar, ushort, uint, ulong, ulonglong, …, ptrdiff
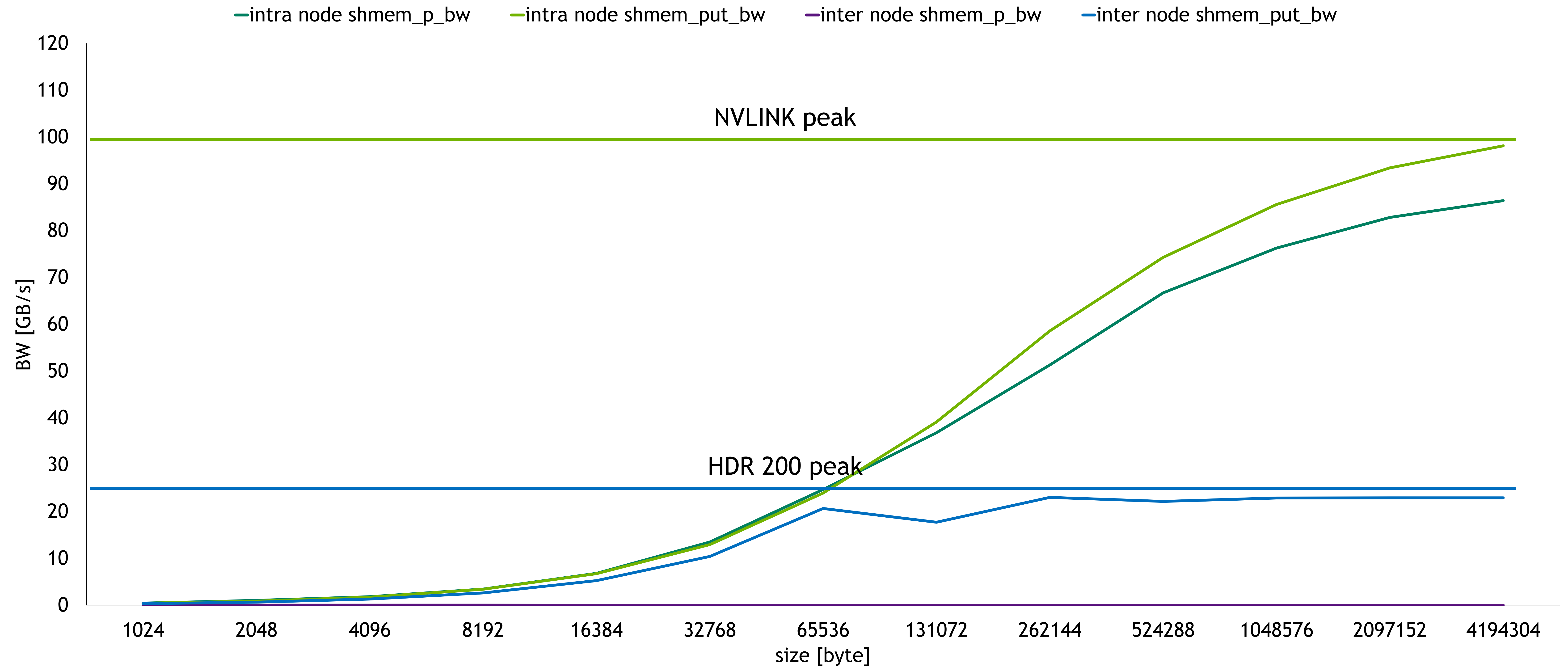(see: https://docs.nvidia.com/hpc-sdk/nvshmem/api/docs/gen/api/rma.html#stdrmatypes)

# NVSHMEM PERFTESTS
## shmem_p_bw and shmem_put_bw on JUWELS Booster – NVIDIA A100 40 GB



—intra node shmem_p_bw  —intra node shmem_put_bw  —inter node shmem_p_bw  —inter node shmem_put_bw

*shmem_p_bw|shmem_put_bw from NVSHMEM perftest*
*JUWELS Booster: https://apps.fz-juelich.de/jsc/hps/juwels/booster-overview.html*

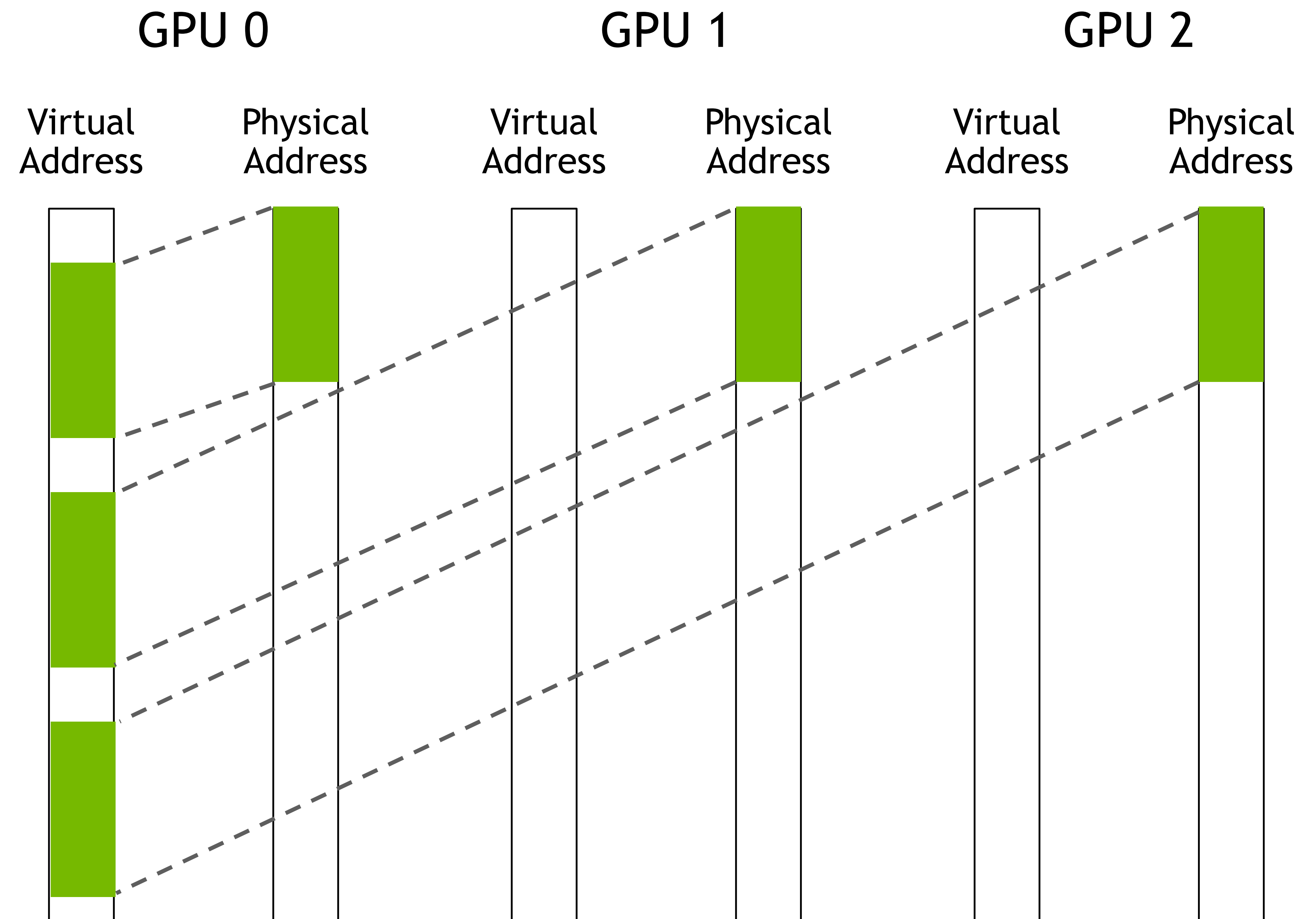Slide 14

NVSHMEM PERFTESTS

shmem_p_bw and shmem_put_bw on JUWELS Booster – NVIDIA A100 40 GB

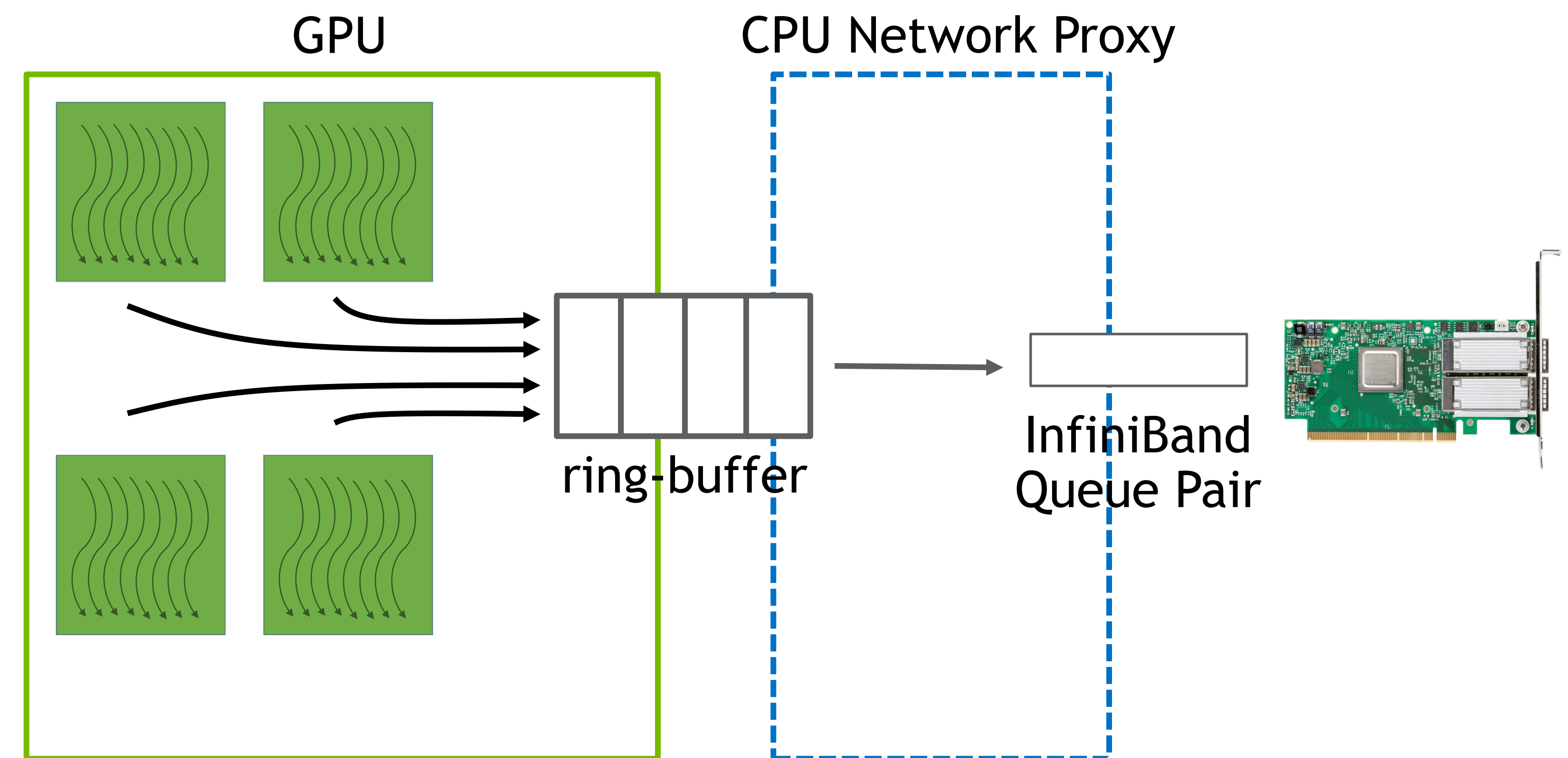Slide 15

# OPTIMIZED INTRA-NODE COMMUNICATION

- Supported on NVLink and PCI-E

- Use CUDA IPC or `cuMem*` API to map symmetric memory of intra-node PEs into virtual address space

- `nvshmem_[put|get]` on device -> load/store

- `nvshmem_[put|get]_on_stream` -> cudaMemcpyAsync



GPU 0    GPU 1    GPU 2

Virtual Address    Physical Address    Virtual Address    Physical Address    Virtual Address    Physical Address
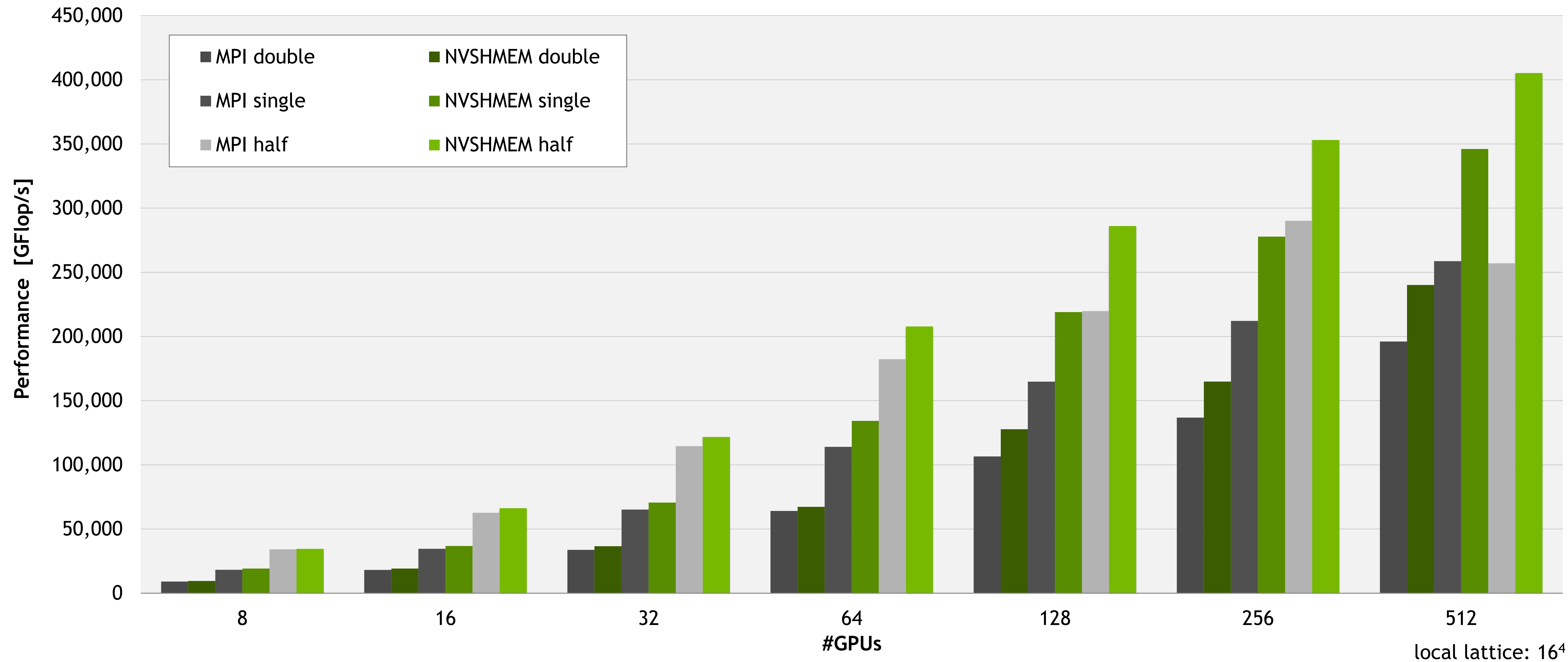
# OPTIMIZED INTER-NODE COMMUNICATION

- NVSHMEM supports inter-node communication over InfiniBand, RoCE, and UCX (experimental)

- Using GPUDirect RDMA (data plane)

- Reverse offloads network transfers from GPU to the CPU (control plane)

- Ring buffer implementation avoids memory fences when interacting with CPU network proxy



GPU

CPU Network Proxy

ring-buffer

InfiniBand Queue Pair

# QUDA STRONG SCALING ON SELENE
## Lattice Quantum ChromoDynamics



local lattice: $16^4$

up to 1.6x Speedup (512 GPUs half precision)

# SUMMARY AND MORE INFORMATION

- Device-initiate communication enables:
    - fine grained communication and computation overlap with sometimes less coding effort
    - kernel fusion not possible with host initiate communication models like MPI and NCCL

- For good intranode device-initiated communication performance it is necessary to aggregate larger messages (`nvshmemx_TYPENAME_put_nbi_block`)

- S32515 NVSHMEM: GPU-Integrated Communication for NVIDIA GPU Clusters: https://www.nvidia.com/en-us/on-demand/session/gtcspring21-s32515/

- S21673 Overcoming Latency Barriers: Strong Scaling HPC Applications with NVSHMEM: https://www.nvidia.com/en-us/on-demand/session/gtcsj20-s21673/

- https://developer.nvidia.com/blog/scaling-scientific-computing-with-nvshmem/

- https://developer.nvidia.com/blog/accelerating-nvshmem-2-0-team-based-collectives-using-nccl/