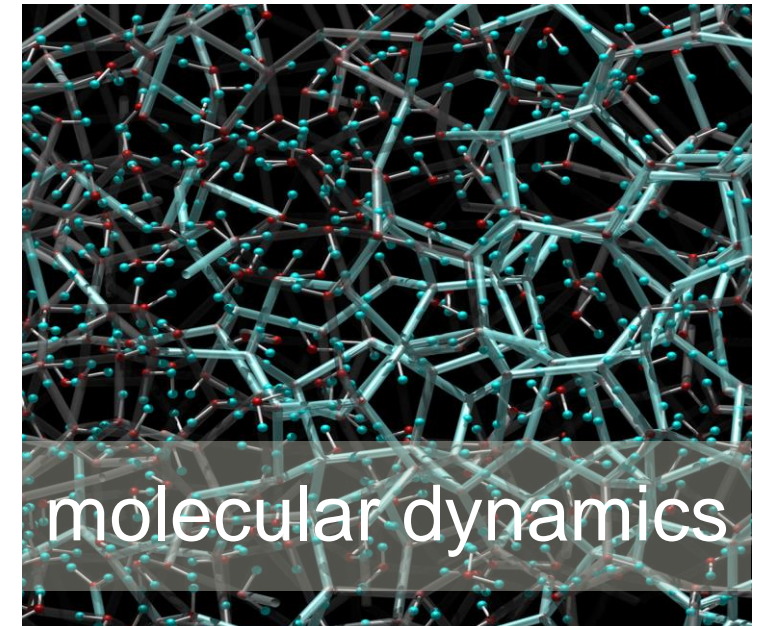# dCUDA: Distributed GPU Computing with Hardware Overlap
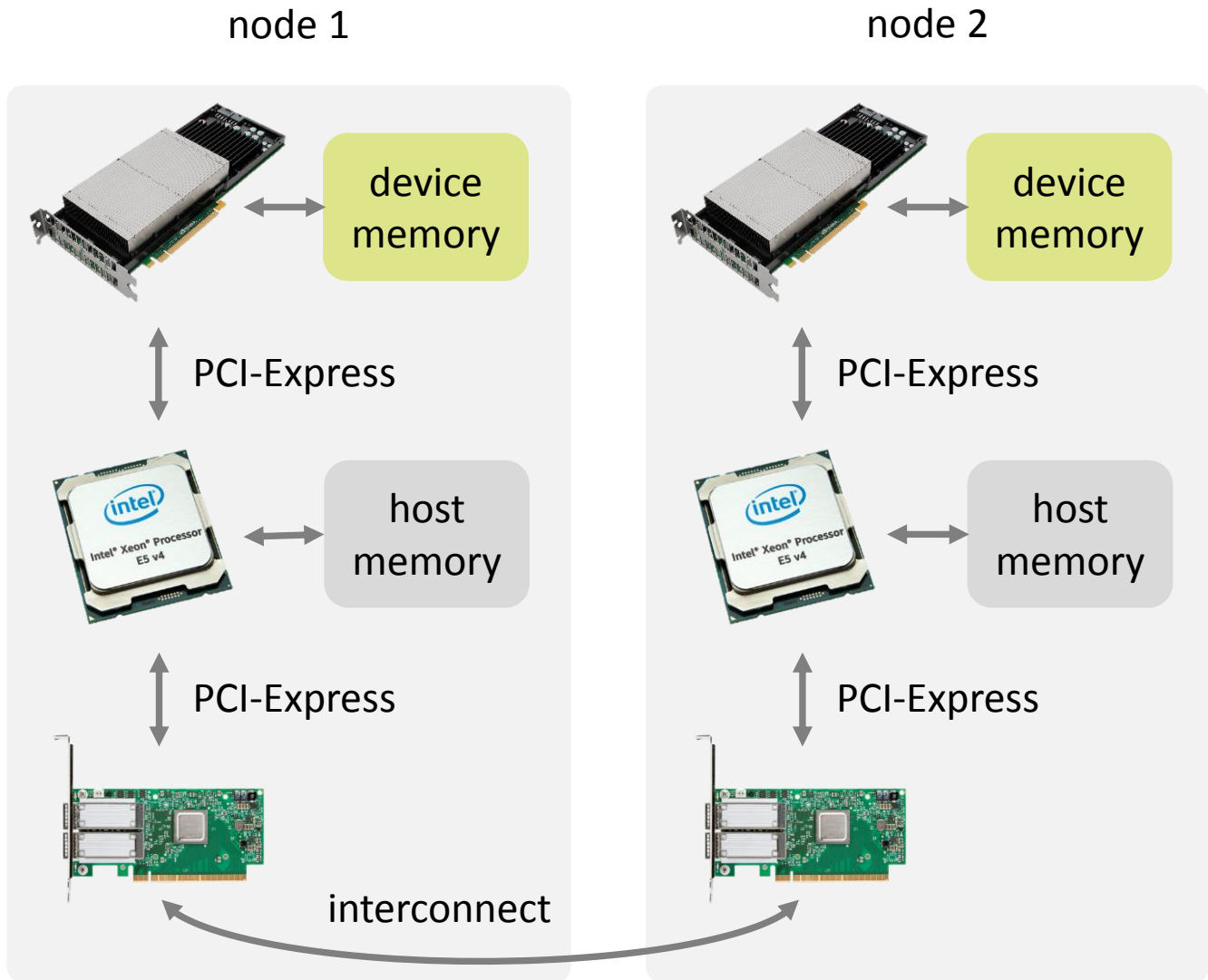
Tobias Gysi, Jeremia Bär, Lukas Kuster, and **Torsten Hoefler**

pasc
Platform for Advanced Scientific Computing

√UK·CUS
Swiss university conference

ETH·RAT

CSCS

# GPU computing gained a lot of popularity in various application domains

weather & climate

machine learning

molecular dynamics

# GPU cluster programming using MPI and CUDA



node 1

node 2

code

```
// run compute kernel
__global__
void mykernel( … ) { }
```
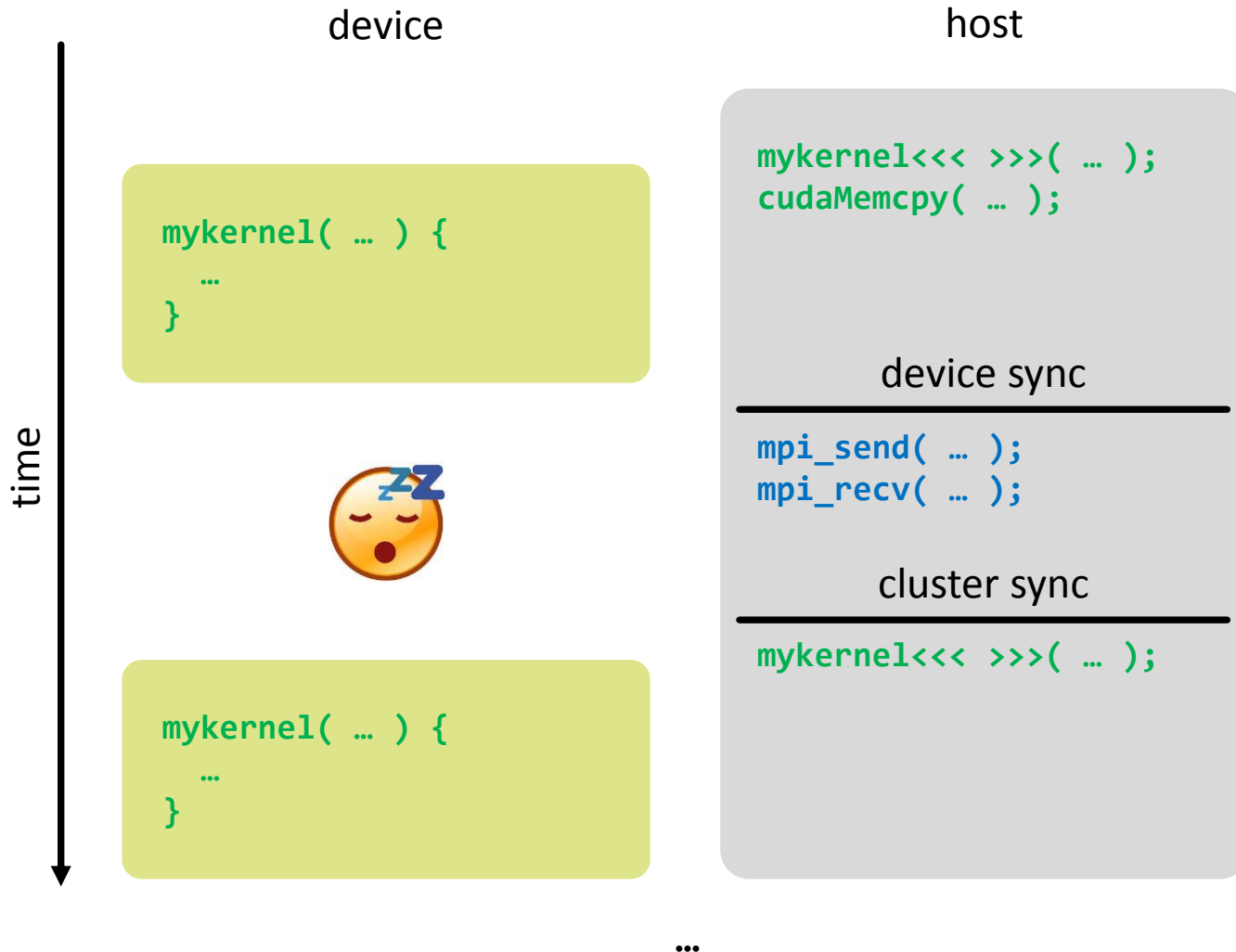
```
// launch compute kernel
mykernel<<<64,128>>>( … );

// on-node data movement
cudaMemcpy(
    psize, &size,
    sizeof(int),
    cudaMemcpyDeviceToHost);

// inter-node data movement
mpi_send(
    pdata, size,
    MPI_FLOAT, … );
mpi_recv(
    pdata, size,
    MPI_FLOAT, … );
```

device memory

host memory

PCI-Express

interconnect

# Disadvantages of the MPI-CUDA approach

device

host

time

```
mykernel( … ) {
    …
}
```

```
mykernel<<< >>>( … );
cudaMemcpy( … );
```

device sync

```
mpi_send( … );
mpi_recv( … );
```

cluster sync

```
mykernel<<< >>>( … );
```

```
mykernel( … ) {
    …
}
```

…

**complexity**
- two programming models
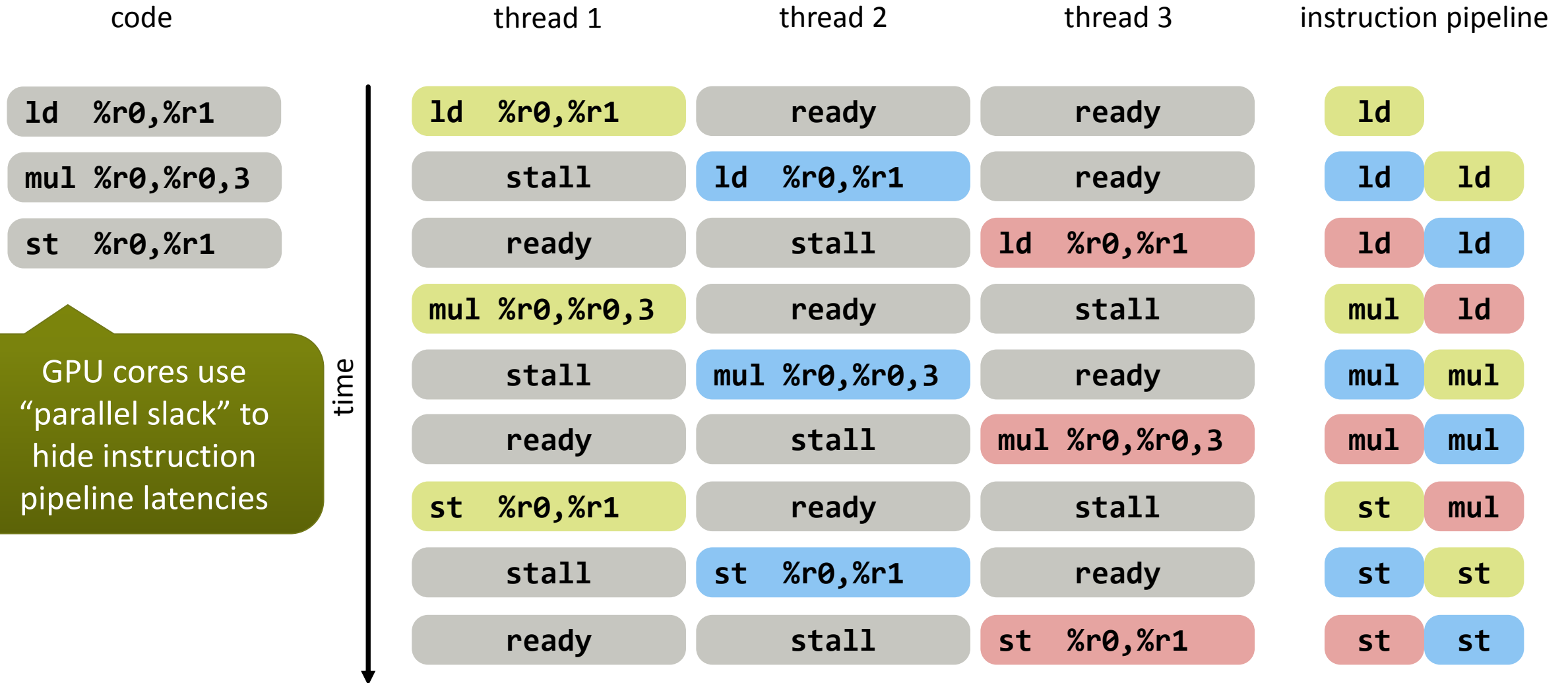- duplicated functionality

copy          sync

**performance**
- encourages sequential execution
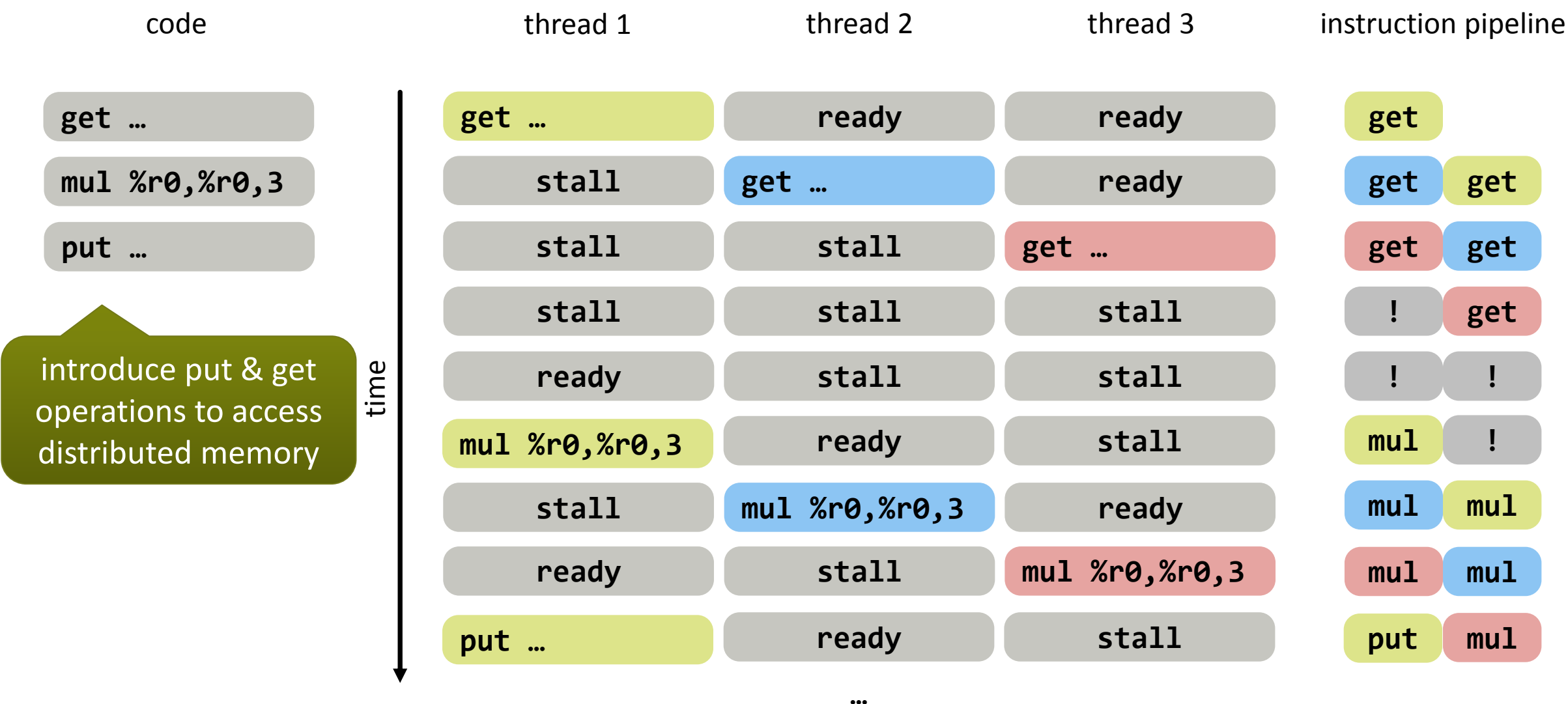- low utilization of the costly hardware

# Achieve high resource utilization using oversubscription & hardware threads

| code | thread 1 | thread 2 | thread 3 | instruction pipeline |
|------|----------|----------|----------|----------------------|
| `ld  %r0,%r1` | `ld  %r0,%r1` | ready | ready | `ld` |
| `mul %r0,%r0,3` | stall | `ld  %r0,%r1` | ready | `ld` `ld` |
| `st  %r0,%r1` | ready | stall | `ld  %r0,%r1` | `ld` `ld` |
| | `mul %r0,%r0,3` | ready | stall | `mul` `ld` |
| | stall | `mul %r0,%r0,3` | ready | `mul` `mul` |
| | ready | stall | `mul %r0,%r0,3` | `mul` `mul` |
| | `st  %r0,%r1` | ready | stall | `st` `mul` |
| | stall | `st  %r0,%r1` | ready | `st` `st` |
| | ready | stall | `st  %r0,%r1` | `st` `st` |

time

GPU cores use "parallel slack" to hide instruction pipeline latencies

...

# Use oversubscription & hardware threads to hide remote memory latencies

| code | thread 1 | thread 2 | thread 3 | instruction pipeline |
|------|----------|----------|----------|----------------------|
| get … | get … | ready | ready | get |
| mul %r0,%r0,3 | stall | get … | ready | get / get |
| put … | stall | stall | get … | get / get |
| | stall | stall | stall | ! / get |
| | ready | stall | stall | ! / ! |
| | mul %r0,%r0,3 | ready | stall | mul / ! |
| | stall | mul %r0,%r0,3 | ready | mul / mul |
| | ready | stall | mul %r0,%r0,3 | mul / mul |
| | put … | ready | stall | put / mul |

introduce put & get operations to access distributed memory

time

…

# How much "parallel slack" is necessary to fully utilize the interconnect?

Little's law

$$concurrency = latency * throughput$$

**device memory**

| latency | ~1µs |
|---|---|
| bandwidth | 200GB/s |
| concurrency | 200kB |
| #threads | ~12000 >> |

# dCUDA (distributed CUDA) extends CUDA with MPI-3 RMA and notifications

```
for (int i = 0; i < steps; ++i) {
  for (int idx = from; idx < to; idx += jstride)
    out[idx] = -4.0 * in[idx] +
        in[idx + 1] + in[idx - 1] +
        in[idx + jstride] + in[idx - jstride];

  if (lsend)
    dcuda_put_notify(ctx, wout, rank - 1,
        len + jstride, jstride, &out[jstride], tag);
  if (rsend)
    dcuda_put_notify(ctx, wout, rank + 1,
        0, jstride, &out[len], tag);

  dcuda_wait_notifications(ctx, wout,
        tag, lsend + rsend);

  swap(in, out);
  swap(win, wout);
}
```

**computation**

**communication**

- iterative stencil kernel
- thread specific idx



- map ranks to blocks
- device-side put/get operations
- notifications for synchronization
- shared and distributed memory

T. Gysi, J. Baer, TH: dCUDA: Hardware Supported Overlap of Computation and Communication, SC16
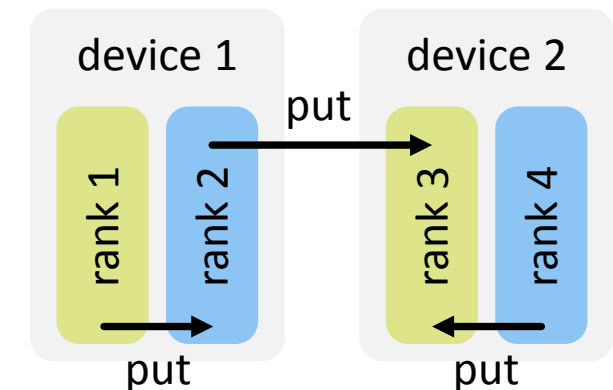
# Advantages of the dCUDA approach



**performance**
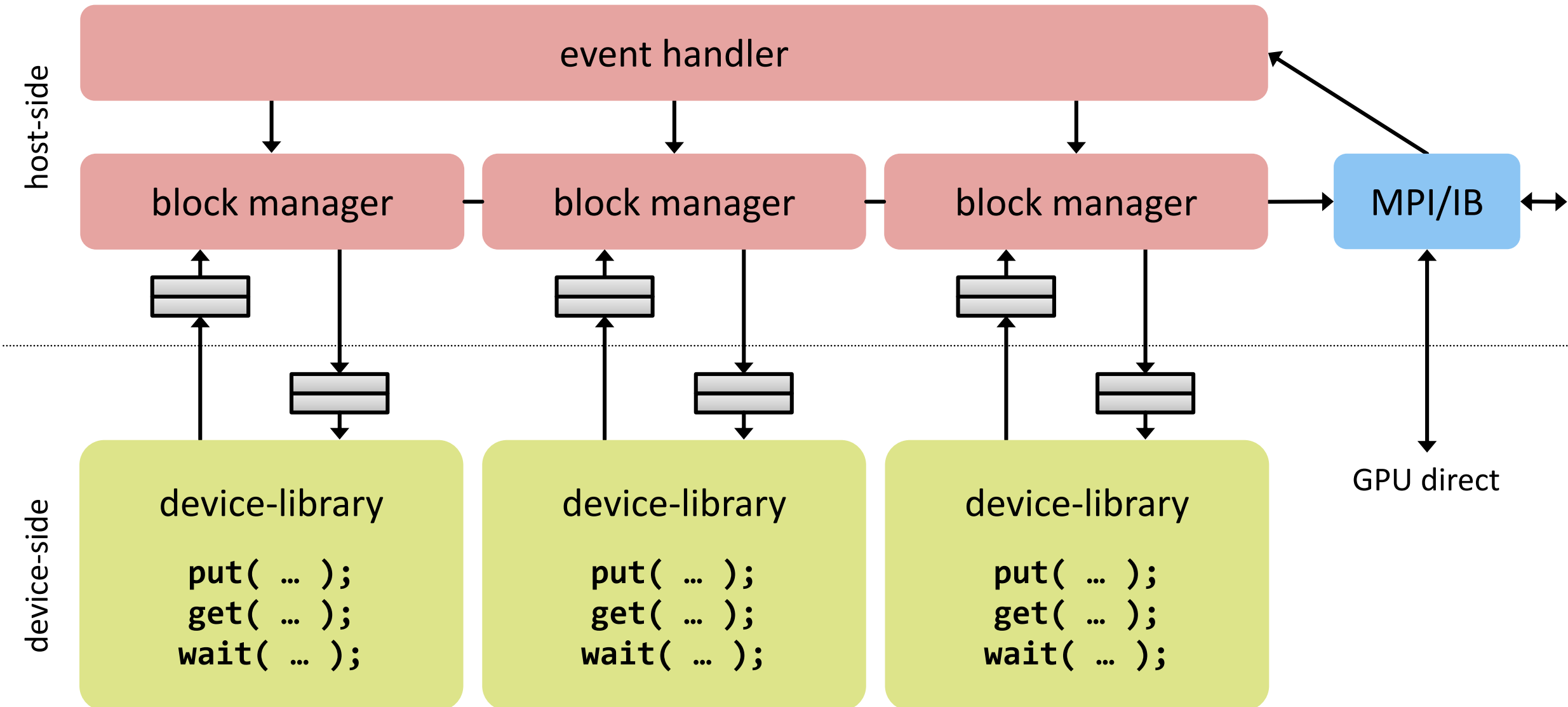- avoid device synchronization
- latency hiding at cluster scale

**complexity**
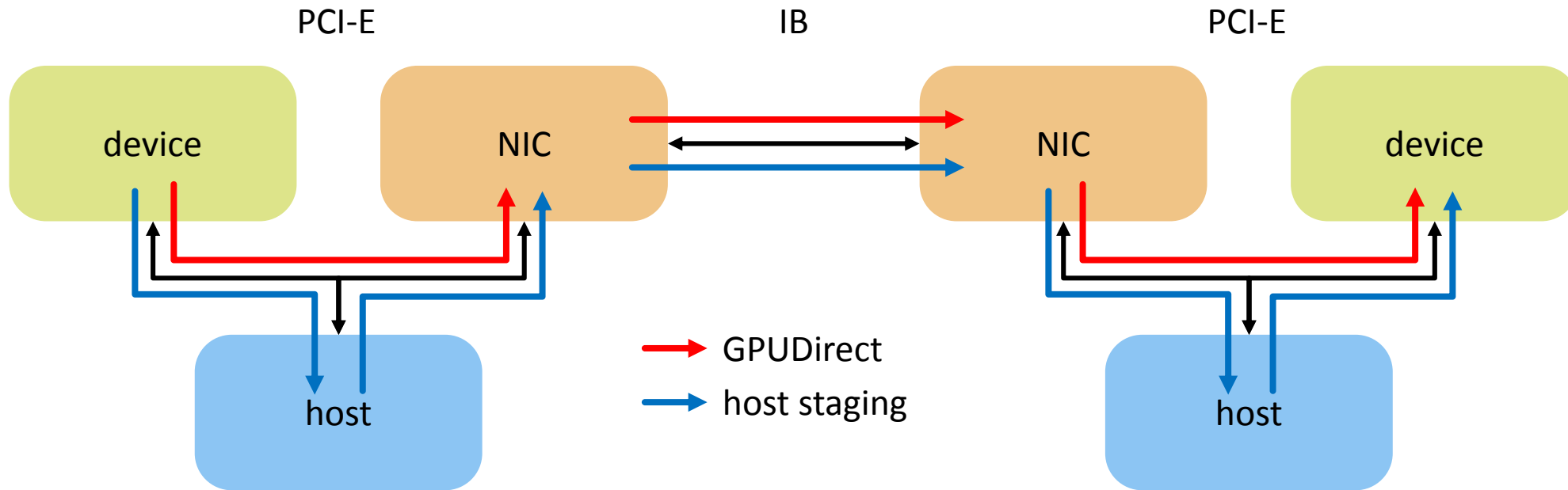- unified programming model
- one communication mechanism

T. Gysi, J. Baer, TH: dCUDA: Hardware Supported Overlap of Computation and Communication, SC16

# Implementation of the dCUDA runtime system

**host-side**

| event handler |

| block manager | block manager | block manager | MPI/IB |

GPU direct

**device-side**

**device-library**

```
put( … );
get( … );
wait( … );
```

**device-library**

```
put( … );
get( … );
wait( … );
```

**device-library**

```
put( … );
get( … );
wait( … );
```

T. Gysi, J. Baer, TH: dCUDA: Hardware Supported Overlap of Computation and Communication, SC16

# GPUDirect provides the NIC with direct device memory access



PCI-E                    IB                    PCI-E

GPUDirect

host staging

**idea**
- avoid copy to host memory
- host-side control

**performance**
- lower latency
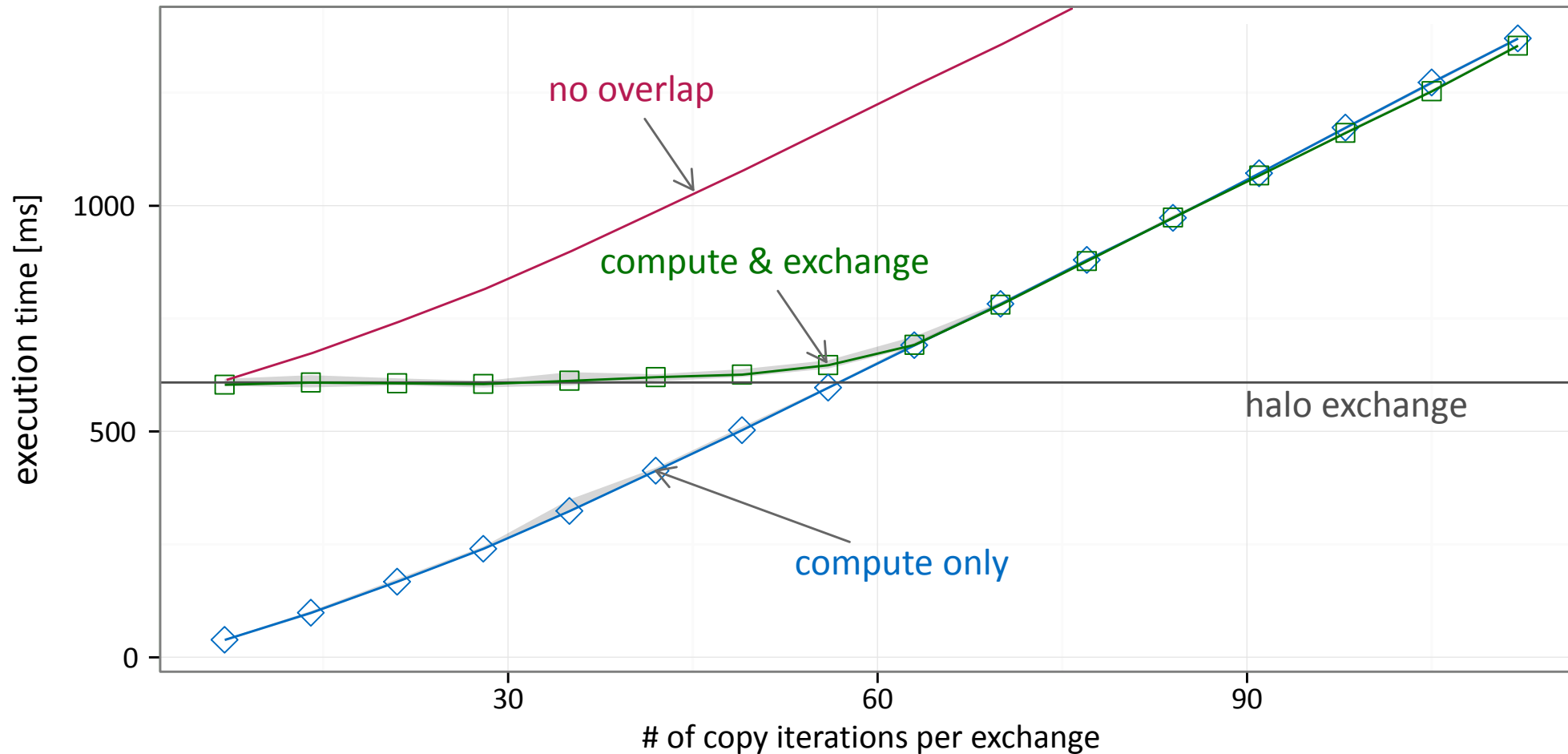- bandwidth penalty on Greina (2.7GB/s instead of 7.2GB/s)

# System latencies of the IB-backend compared to the MPI-backend

benchmarked on Greina (4 Broadwell nodes with 1x Tesla K80 per node)

| | IB-backend | MPI-backend |
|---|---|---|
| same device (put & notify) [μs] | 2.4 | 2.4 |
| peer device (put & notify) [μs] | 6.7 | 23.7 |
| remote device (put & notify) [μs] | **6.9** | **12.2** |
| same device (notify) [μs] | 1.9 | 1.9 |
| peer device (notify) [μs] | 3.4 | 5.0 |
| remote device (notify) [μs] | **3.6** | **5.4** |

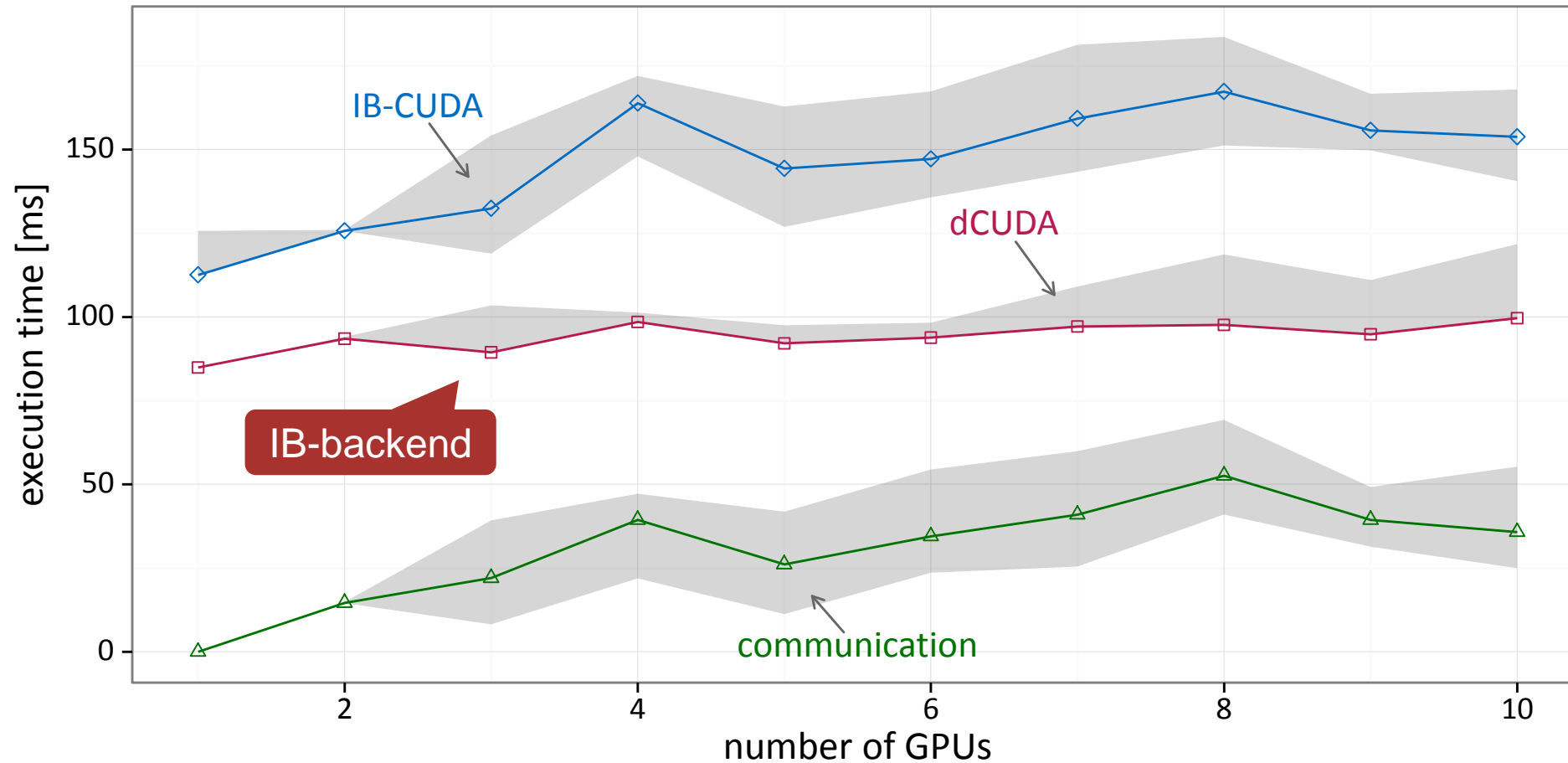T. Gysi, J. Baer, TH: dCUDA: Hardware Supported Overlap of Computation and Communication, SC16

# Overlap of a copy kernel with halo exchange communication

benchmarked on Greina (8 Haswell nodes with 1x Tesla K80 per node)



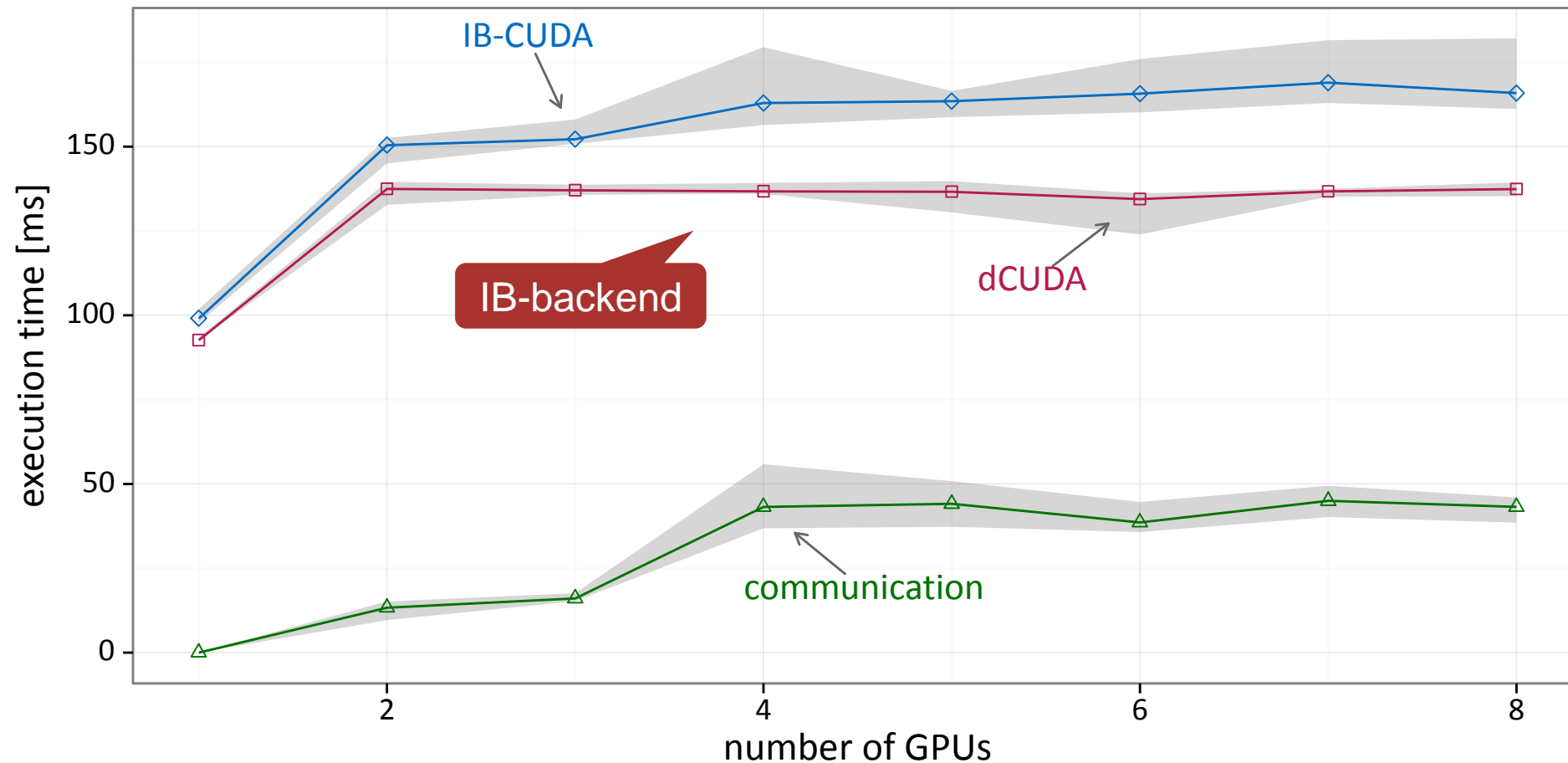T. Gysi, J. Baer, TH: dCUDA: Hardware Supported Overlap of Computation and Communication, SC16

# Weak scaling of IB-CUDA and dCUDA for a particle simulation

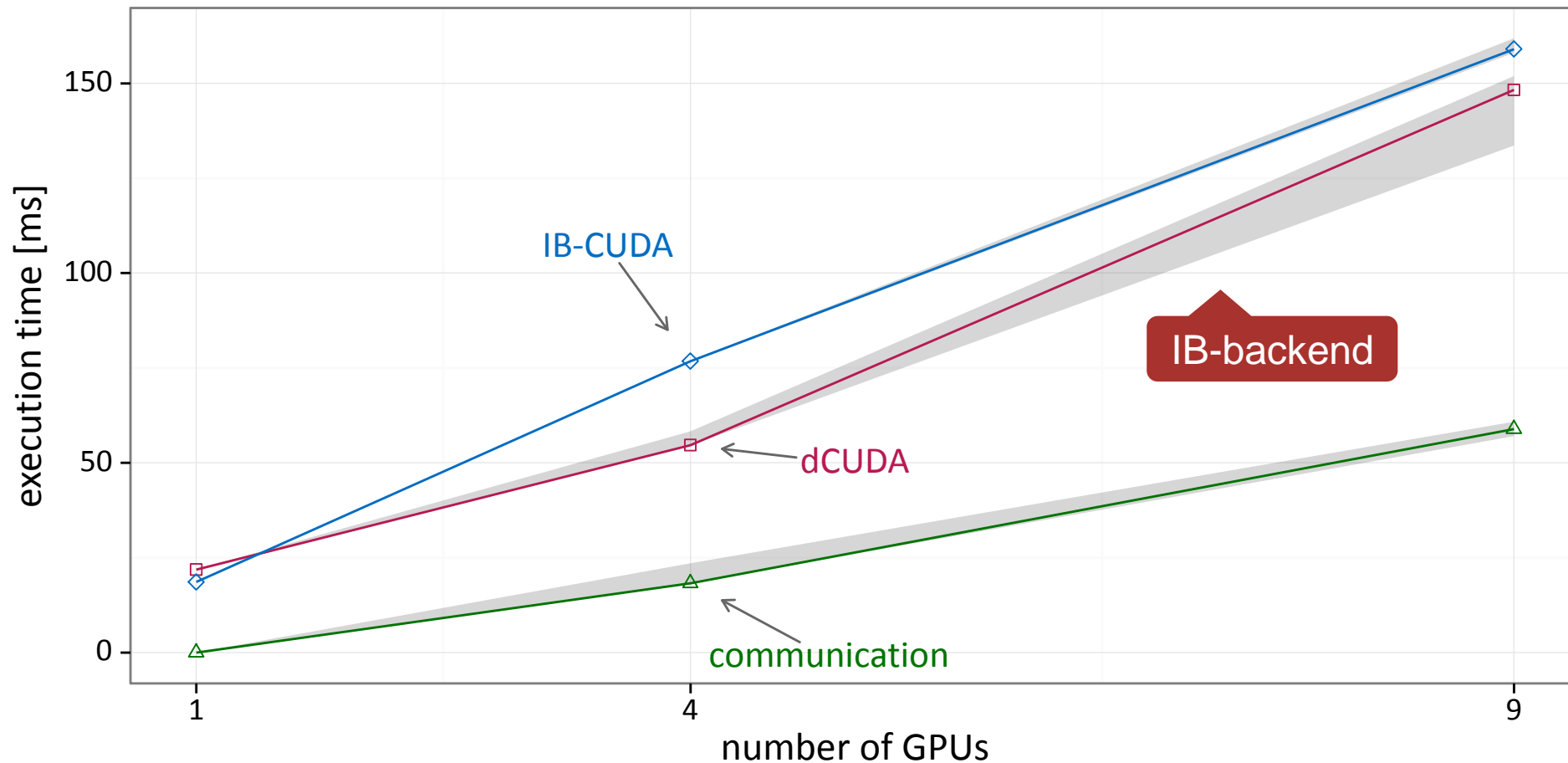benchmarked on Greina (4 Broadwell nodes with 1x Tesla K80 per node)

# Weak scaling of IB-CUDA and dCUDA for a stencil program

benchmarked on Greina (4 Broadwell nodes with 1x Tesla K80 per node)

# Weak scaling of IB-CUDA and dCUDA for sparse-matrix vector multiplication
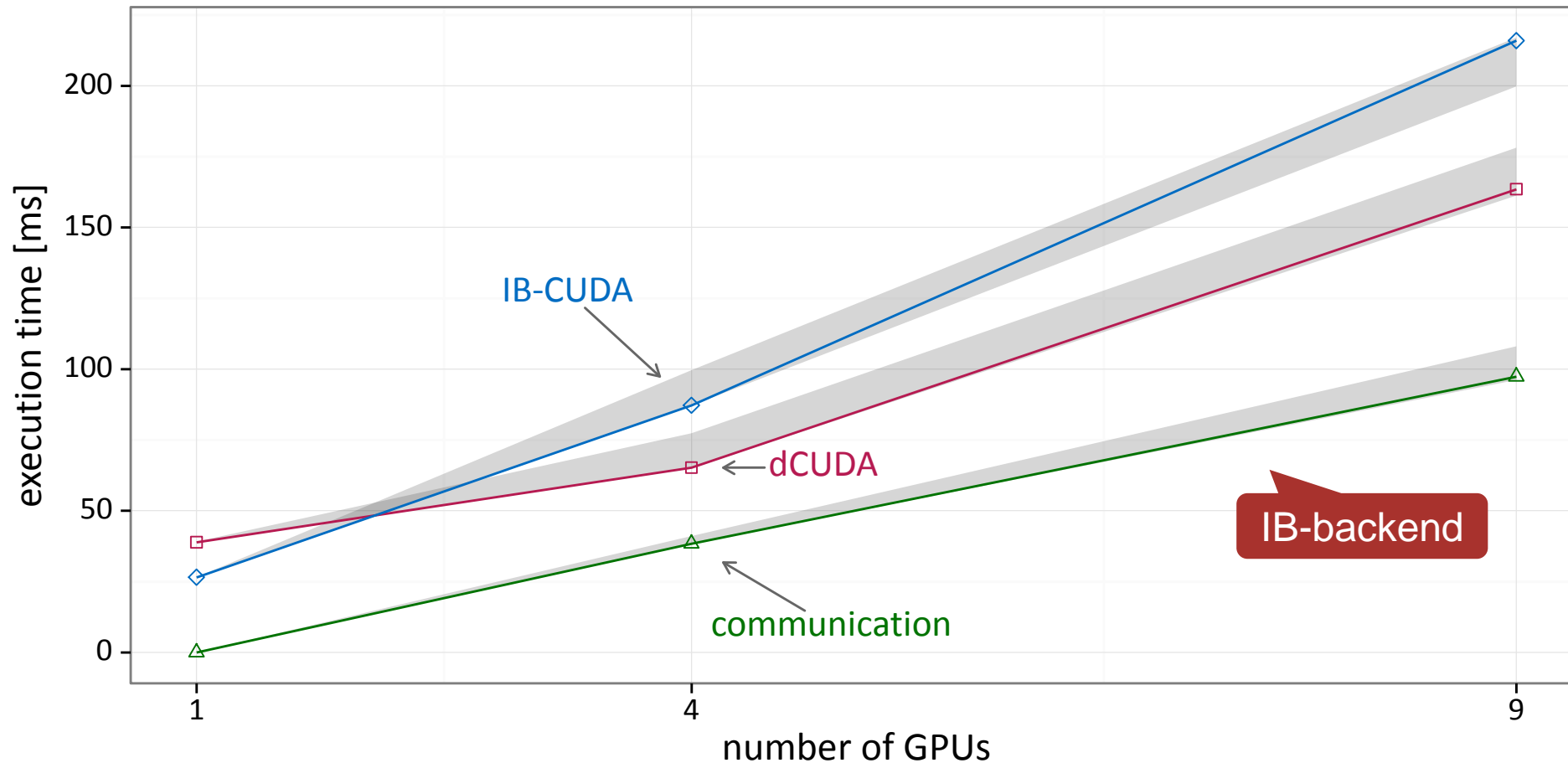
benchmarked on Greina (4 Broadwell nodes with 1x Tesla K80 per node)

# Weak scaling of IB-CUDA and dCUDA for power iterations

benchmarked on Greina (4 Broadwell nodes with 1x Tesla K80 per node)

# Conclusions

- unified programming model for GPU clusters
    - device-side remote memory access operations with notifications
    - transparent support of shared and distributed memory
- extend the latency hiding technique of CUDA to the full cluster
    - inter-node communication without device synchronization
    - use oversubscription & hardware threads to hide remote memory latencies
- automatic overlap of computation and communication
    - synthetic benchmarks demonstrate perfect overlap
    - example applications demonstrate the applicability to real codes
- https://spcl.inf.ethz.ch/Research/Parallel_Programming/dCUDA/

Platform for Advanced Scientific Computing    Swiss university conference    ETH-RAT    CSCS

T. Gysi, J. Baer, TH: dCUDA: Hardware Supported Overlap of Computation and Communication, SC16