

Extending MPI to Accelerators*

Jeff A. Stuart, John D. Owens
University of California, Davis

cpunerd@gmail.com, jowens@ece.ucdavis.edu

Pavan Balaji
Argonne National Laboratories
balaji@mcs.anl.gov

Outline

- Motivation
- Previous Work
- Proposal
- Challenges

Motivation

- HPC no longer (just) CPU
- GPUs Have Problems
 - Slave Device
 - No system calls

Previous Work

- Three Main Works
 - cudaMPI
 - GAMPI
 - DCGN

Previous Work

- cudaMPI
 - Handles buffer movement
 - No ranks for GPUs

Previous Work

- GAMPI
 - GPUs have ranks*
 - More communicators
 - Handles buffer movement

Previous Work

- DCGN
 - GPUs have ranks
 - GPUs source/sink communication*
 - Doesn't implement standard MPI

Proposal

- Several Ideas
 - No Ranks for GPUs
 - Multiple Ranks per GPU Context
 - One Rank per GPU Context
 - New MPI Function(s) to Spawn Kernels

Proposal

- No Ranks for GPUs
 - The way things work right now
 - No changes necessary to MPI

Proposal

- Multiple Ranks Per Accelerator Context
 - Ranks exist for lifetime of application
 - # of ranks chosen at runtime by user
 - Modifications to MPI
 - Bind GPU threads to rank/MPI functions take source rank
 - Host must listen for requests
- Extra threads on CPU (one for each GPU)

Proposal

- One Rank per Accelerator Context
 - Ranks exist for lifetime of Application
 - Mapping of Processes:Contexts?
 - Can CPU Processes use MPI communication?

Proposal

- New MPI Function(s) to Spawn Kernels
 - New communicators and ranks after every spawn
 - Cleaned up after all kernels finish
 - Intercommunicator(s) available upon request

Challenges

- Threads vs Processes
- Extra Communicators?
- Collectives
- Source/Sink Communication

Looking Forward

- GPU-Direct is good
- GPU-Direct 2 is great
- We want GPU-Direct 3 to
 - Let GPU source/sink
 - Use GPU-Direct 2 to interface with NIC
 - Administer MPI ranks without CPU interference

One Last Note

- Graduating with Ph.D. In June 2012
- Resume at <http://jeff.bleugris.com/resume.pdf>