

# Interconnect Bandwidth Heterogeneity on AMD MI250x and Infinity Fabric

Carl Pearson

*Sandia National Laboratories*

Albuquerque, NM, USA

cwpears@sandia.gov

**Abstract**—Demand for low-latency and high-bandwidth data transfer between GPUs has driven the development of multi-GPU nodes. Physical constraints on the manufacture and integration of such systems has yielded heterogeneous intra-node interconnects, where not all devices are connected equally. The next generation of supercomputing platforms are expected to feature AMD CPUs and GPUs. This work characterizes the extent to which interconnect heterogeneity is visible through GPU programming APIs on a system with four AMD MI250x GPUs, and provides several insights for users of such systems.

## I. INTRODUCTION

Multi-GPU nodes are a common feature of high-performance computing systems due to their density and power-efficiency. For suitable workloads, GPUs offer more floating-point operations per watt, and more operations per second available within a single node. Physical constraints place limits on the number and bandwidth of chip-to-chip interconnects. Therefore, multi-CPU and multi-GPU nodes may be organized with heterogeneous interconnects between different pairs of components to maximize the amount of compute available in the node. Systems with AMD x86 CPUs and AMD GPUs are expected to be forerunners in the upcoming generation of supercomputing systems. Frontier [1] is an example of such a system; a familiar design, with a single CPU connected to four GPUs by AMD’s proprietary Infinity Fabric interconnect.

This work contributes

- the first published relationship between communication primitives and achieved bandwidth on AMD’s CDNA 2 GPUs and Infinity Fabric 3 interconnect.

Furthermore, it makes the following observations from the results:

- like previous multi-GPU systems, interconnect heterogeneity manifests at the HIP API level as significant

Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy’s National Nuclear Security Administration under contract DE-NA-0003525. This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration. This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725. The author would like to thank James Elliot and Simon Garcia de Gonzalo of Sandia National Labs for their feedback and guidance.

bandwidth differences between HIP devices depending on which devices are participating,

- unlike previous multi-CPU, multi-GPU systems, NUMA effects were not observed between CPU and GPU,
- if GPUs are idle during communication, execution resources can be used to achieve higher bandwidth.

The rest of this paper is organized as follows. Section II discusses the measurement methods, Sec. III discusses the results, Sec III-F discusses related work, and Sec. III-G concludes.

## II. METHODOLOGY

### A. Evaluation System

AMD provides the ROCm software platform [2] for GPU-accelerated systems. It includes C++ language extensions, kernel language, and runtime APIs (together known as HIP) necessary to interact with AMD GPUs, as well as higher-level libraries. This work uses the HIP runtime API and kernel language to evaluate point-to-point CPU/GPU communication bandwidth on a flagship HPC node: Crusher [3] at the Oak Ridge Leadership Computing Facility (OLCF), which has nodes identical to Frontier.

Feature	Description
CPU	AMD EPYC 7A53
GPU	4x AMD MI250x (2x GCD)
GCD	4x2 AMD CDNA2
CPU-GCD	Infinity Fabric 72+72 GB/s
Intra-GPU (“quad”)	Infinity Fabric 200+200 GB/s
Inter-GPU (“dual”)	Infinity Fabric 100+100 GB/s
Inter-GPU (“single”)	Infinity Fabric 50+50 GB/s
ROCm	5.4.0
GPU Driver	5.16.9.22.20.7582
OS	SUSE Linux Enterprise Server 15 SP3
Kernel	Linux 5.3.18

TABLE I: OLCF Crusher node summary. Bandwidths are given as the sum of each direction. Figure 1 summarizes the node topology. “Quad”, “dual”, and “single” refer to the number of connections drawn in that figure.

Table I and Figure 1 summarize the topology of the node. Each node has a single AMD EPYC 7A53 CPU with four AMD MI250x GPUs (the “host” and “devices” respectively). Each MI250x has two graphics compute dies (GCDs), each of which can be addressed as a unique HIP device. CPUs and GCDs are connected by different kinds of Infinity Fabric links [4].

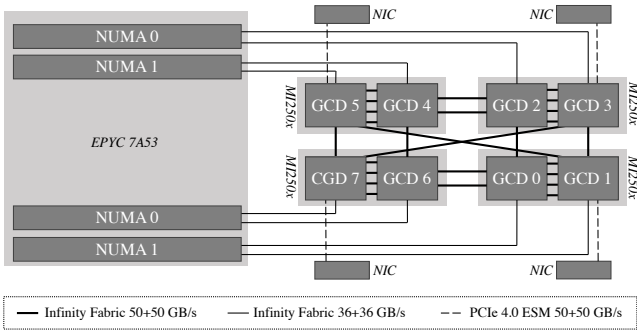


Fig. 1: Crusher node block-diagram, adapted from [3]. Each of the four AMD MI250x GPUs has two graphics compute dies (GCDs), each of which is an addressable HIP device acting as an individually-programmable GPU. The system features PCIe 4.0, as well as bidirectional 50 + 50 GB/s and 36 + 36 GB/s Infinity Fabric interconnects.

The in-package infinity fabric offers 200+200 GB/s bidirectional bandwidth to neighboring GCD (“quad”). There are also 8 lanes of inter-package Infinity Fabric, for 400+400 GB/s total. On the evaluation platform, this is allocated in two 100+100 GB/s (“dual”) and one 50+50 GB/s (“single”) connection to other GCDs and one 36+36 GB/s coherent connection to an L3 slice on the CPU. There is also a 50+50 GB/s PCIe 4.0 Extended Speed Mode connection to NIC not investigated in this work.

### B. Allocation Types

Table II summarizes the buffer types and transfer methods in this work. All allocations are coarse-grained, i.e. they do not support any coherency during GPU kernel execution.

Device allocations are produced with `hipMalloc`. They may be used for an explicit transfer, or mapped into another HIP device using `hipPeerAccessEnable` so a kernel executing on that device can access the data implicitly.

Host pinned allocations are produced by `hipHostMalloc` with `hipHipHostMallocNumaUser` and `hipHipHostMallocNonCoherent` flags. The buffer can be used in an explicit transfer, or accessed implicitly by a GPU kernel with `hipHostGetDevicePointer`. The data in this allocation may not be paged by the host.

Host pageable allocations are produced with the system allocator (e.g., `malloc`). This buffer can be used in an explicit transfer; internally HIP runtime will stage the data through a pinned buffer so the GPU DMA engine can read it from physical memory.

Managed allocations are produced with `hipMallocManaged` and `hipMemAdviseSetCoarseGrain`. This buffer may be transparently accessed by the GPU and the CPU with various degrees of coherency.

### C. Transfer Methods

*Explicit* transfers use `hipMemcpyAsync`. A buffer is created on the source and destination, and `hipMemcpyAsync` is invoked to move data from source to destination.

*Implicit* transfers use load and store operations in a GPU kernel to move data from source to destination. For *implicit mapped* transfers between host and device, a buffer is created on the host, and the device writes to or reads from that buffer. For device-device transfers, the buffer is created on the destination device, and the source writes to it. For *implicit managed* transfers, a single managed allocation is produced. It is prefetched to the source device, and then modified from the destination device. `HSA_XNACK=1` is set in the environment, causing pages to be migrated from the source to destination device. The `cpu_write` function uses an OpenMP-parallel loop to write 64-bit elements to the buffer with sequential values equal to their index plus a fixed constant. The `gpu_write` kernel does the same using a large HIP grid, where threads make coalesced accesses to the buffer. The `gpu_read` kernel instead reads those values.

*Prefetch* transfers use `hipMemPrefetchAsync` on a managed buffer so the data is resident on a particular device. The data will then be local to that device if that device makes future accesses.

Host Buffer	Device Buffer	Transfer	Direction
<code>malloc</code>	<code>hipMalloc</code>	explicit	H2D, D2H
<code>hipHostMalloc</code>	<code>hipMalloc</code>	explicit	H2D, D2H
<code>hipHostMalloc</code>	–	implicit	H2D, D2H
–	<code>hipMalloc</code>	explicit	D2D
–	<code>hipMalloc</code>	implicit	D2D
<code>hipMallocManaged</code>	–	implicit	H2D, D2H, D2D
<code>hipMallocManaged</code>	–	prefetch	H2D, D2H, D2D

TABLE II: Combinations of buffer types and transfer methods evaluated in this work. “Explicit” transfers use `hipMemcpyAsync`. “Implicit” transfers use load/store operations on a buffer pointer in a GPU kernel. “D2H,” “H2D,” and “D2D” mean device-to-host, host-to-device, and device-to-device respectively.

### D. Measurement

Measurements are implemented in the `Comm|Scope` package at [github.com/c3sr/comm\\_scope](https://github.com/c3sr/comm_scope). The Google Benchmark support library [5] is used to drive the various benchmarks. It chooses the number of measurements iterations such that the operation in question executes for at least one second, at least once, and fewer than one billion times. For these benchmarks, the fastest (GPU-to-GPU implicit writes) were invoked  $\approx 59000$  times, and the slowest (prefetches) twice.

In the setup phase, NUMA affinity for host allocations is enforced as necessary. `hipDeviceReset` is used to reset the devices to discard any state that may have accumulated from past benchmark runs. `hipSetDevice` controls the active device(s) while necessary buffers are created and filled to ensure a physical memory mapping.

During the benchmark iterations, the benchmark state is reset by a combination of device affinities, cache flushes, memory fills, and prefetches to get the buffers to a known state. For asynchronous operations, the start HIP event is recorded, the operation is invoked in the default stream, and the stop event is recorded. For synchronous operations, an `std::chrono::high_precision` clock is recorded, the operation is invoked, and the clock is recorded again to measure the wall time. This process is repeated until the

Google Benchmark harness’ conditions are satisfied. During teardown, all resources and NUMA bindings are released.

### III. RESULTS

Figures 2a, 2b, and 2c show the measured GCD-to-GCD bandwidth for GCDs 0-1 (quad-link, intra-GPU), 0-6 (dual link), and 0-2 (single-link) respectively. Figure 3a and 3b show the measured CPU-to-GCD and GCD-to-CPU bandwidth, respectively. The primary observed effect is that the transfer method has an enormous impact on achieved bandwidth over fixed hardware. Only implicit access by GPU kernels are able to saturate every interconnect in the system.

#### A. Managed Memory Prefetch Bandwidth

The `hipMemPrefetchAsync` hint can be used to prefetch managed data to a specific device in anticipation of future access there. It is several orders of magnitude (up to 1630 $\times$ , 47 $\times$  for 1 GiB transfers) slower than the fastest transfer method. It is too slow to elucidate any interconnect heterogeneity, and will be excluded from discussion in the sections below.

#### B. Transfer Method Effects

The effect of the transfer method is most significant for the fastest interconnects. Table III shows what fraction of the theoretical peak bandwidth is achieved for each transfer method for different GCD-GCD connections. For faster intra-GPU interconnects, achieved bandwidth varies by 3 $\times$ , while for the slower single-linked GCDs all transfer methods are roughly equivalent, with dual-link GCDs somewhere in-between. Implicit transfers between mapped buffers are always able to saturate the link. This suggests that using the GPU’s computing resources to speed up transfers may be more effective than using slower DMA methods, especially if the GPU is idle during the transfer.

Transfer	Interconnect Type		
	“quad”	“dual”	“single”
explicit	0.25	0.51	0.76
implicit mapped	0.77	0.77	0.78
implicit managed	0.74	0.76	0.76
prefetch managed	0.016	0.032	0.064
<i>Peak GB/s</i>	<i>200</i>	<i>100</i>	<i>50</i>

TABLE III: Fraction of theoretical peak bandwidth for 1 GiB device/device transfers

While the 50 GB/s GCD-GCD link is slow enough to make all transfer methods equal, the even slower CPU-GCD link is still fast enough to reveal the slowing effect of staging pagable transfers through pinned memory. Figures 2 and 3 show that such transfers are 5 $\times$  slower than pinned allocations in the worst case.

#### C. Device/Device Topology Effects

Device-to-device topology effects are especially apparent for the fastest transfer methods. Implicit access to a mapped allocation on a different device is able to utilize  $\approx 75\%$  of the available bandwidth, yielding 153 GB/s within a GPU, 77 GB/s between dual-link GPUs, and 39 GB/s between single-link GPUs. Explicit transfers seem to have a ceiling

of 51 GB/s, so while they are able to saturate the single-link GPU to approximately the same bandwidth as implicit access (38 GB/s), they are only able to generate 51 GB/s across the much higher-bandwidth intra- and inter-GPU connections. This suggests the DMA engine in CDNA2 may only be able to generate 51 GB/s of memory traffic for a given transfer.

#### D. Host/Device Topology Effects

Despite different NUMA regions being linked with different GCDs, there are no observable NUMA effects on transfer bandwidth (i.e., bandwidth is the same regardless of NUMA region or GCD involved). Thus, Figs. 3a and 3b are representative of any host-device pairing. While these results suggest that NUMA affinity of host allocations does not matter, it may become more relevant if multiple transfers are in flight simultaneously.

These experiments cannot reveal whether this homogeneity is a result of the CPU memory subsystem or the GPU interconnects, since the fastest CPU/GPU transfers are slower than the slowest 38 GB/s GPU-GPU transfers, so any GPU-GPU interconnect effects would not be visible.

#### E. Anisotropic Effects

CPU-to-GCD (Fig. 3a, “implicit managed”) transfers are much faster than GCD-to-CPU (3b, “implicit managed”). For the former, the GPU is initiating writes to CPU-resident pages and vis-versa for the latter. This is the only substantial anisotropic transfer observed.

#### F. Related Work

Substantial prior work exists for investigating multi-GPU communication performance. All prior work operates on widespread Nvidia platforms. This work distinguishes itself by doing similar evaluations on AMD platforms.

Spafford, Meredith, and Vetter [6] identify NUMA-related bandwidth effects for CPU/GPU transfers. This work makes similar measurements on a more modern platform and finds that it has no such NUMA effects. Tallent et al. [7] investigate the performance of GPU-GPU data transfers. While they focus on certain collective operations relevant to deep neural networks, they also measure point-to-point transfer time at specific sizes. Li et al. [8] present microbenchmarks of latency, bandwidth, and collective communications for PCIe 3.0, and Nvidia’s NVLink, NVSwitch, and GPUDirect technologies. Pearson et al. [9] do a similar evaluation for Nvidia NVLink interconnects and IBM POWER9 (excluding collective communications). This work does a similar analysis of AMD’s Infinity Fabric interconnect with AMD GPUs. Alan and Ge [10] and Chien, Peng, and Markidis [11] present a characterization of the unified virtual memory system on Nvidia GPUs across microbenchmarks, synthetic kernels, and HPC workloads. This work microbenchmarks achievable bandwidth for AMD’s version of this technology.

#### G. Conclusion and Future Work

This paper presents GPU-GPU communication bandwidth measurements for a flagship AMD high-performance compute

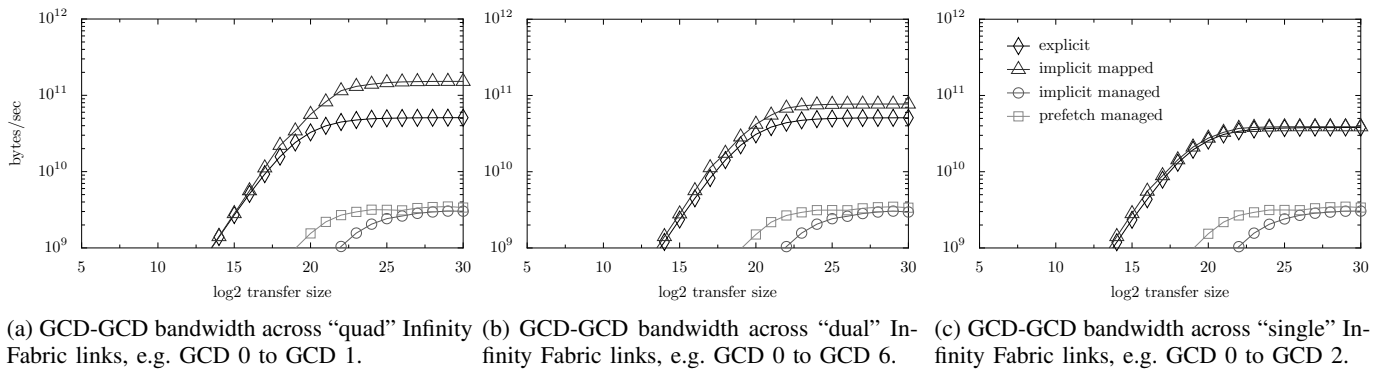


Fig. 2: Unidirectional GCD-to-GCD bandwidth measured on the OLCF Crusher system.

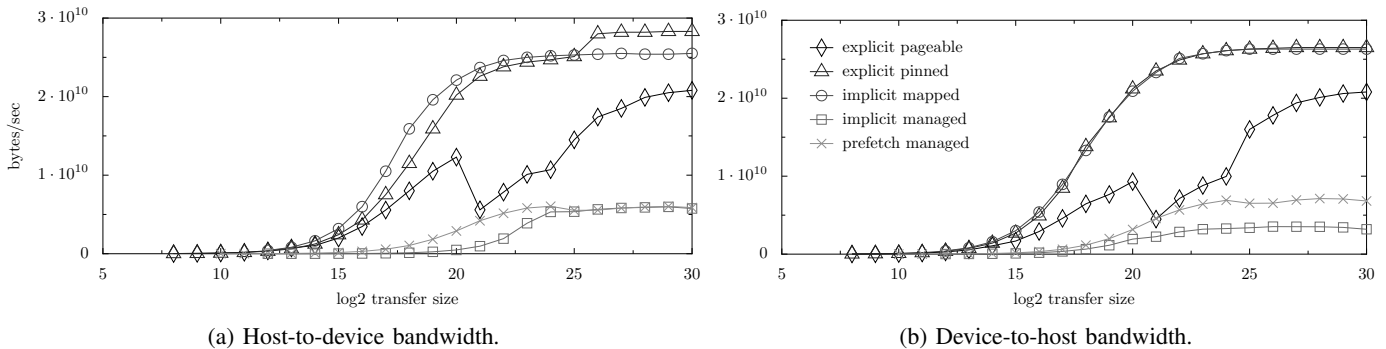


Fig. 3: Unidirectional NUMA-node / GPU bandwidth on OLCF Crusher node. All pairs of NUMA-node and HIP device have the same performance characteristics. Results for NUMA 0 and GCD 0 shown.

node. The capabilities of the intra- and inter-GPU interconnects have outstripped the ability of the system software and DMA engines to saturate them. If the workload requires GPU-to-GPU communication, it may be beneficial to implement it as implicit access between mapped buffers rather than relying on the DMA engine. There were no observable NUMA effects impacting CPU/GPU transfers, which may simplify this dimension of multi-GPU application tuning. The presented methodology is implemented in the CommScope package at [github.com/c3sr/comm\\_scope](https://github.com/c3sr/comm_scope).

Future systems like El Capitan that contain AMD processors are expected to feature integrated CPUs and GPUs, with substantially different performance characteristics despite the same ROCm programming development platform. These systems may feature even higher-bandwidth interconnects and tighter integration, which further emphasize distinctions between transfer methods and impact the performance of implicit accesses between different devices.

The evaluation in this paper is limited to unidirectional point-to-point transfers between CPUs and GPUs using HIP APIs. There are significant additional transfer types of interest to the community, including simultaneous (including bidirectional and collective), GPU-NIC, and off-node.

## REFERENCES

- [1] Oak Ridge Leadership Computing Facility Staff. (2022) Frontier user guide. Oak Ridge Leadership Computing Facility. [Online]. Available: [https://docs.olcf.ornl.gov/systems/frontier\\_user\\_guide.html](https://docs.olcf.ornl.gov/systems/frontier_user_guide.html)
- [2] (2022) ROCm documentation. AMD. [Online]. Available: <https://docs.amd.com/>
- [3] (2022) Crusher quick-start guide. Oak Ridge Leadership Computing Facility. [Online]. Available: [https://docs.olcf.ornl.gov/systems/crusher\\_quick\\_start\\_guide.html](https://docs.olcf.ornl.gov/systems/crusher_quick_start_guide.html)
- [4] (2021) Introducing AMD CDNA 2 architecture. AMD. [Online]. Available: <https://www.amd.com/system/files/documents/amd-cdna2-white-paper.pdf>
- [5] (2022) Google Benchmark Support Library. Google Benchmark authors. [Online]. Available: <https://github.com/google/benchmark>
- [6] K. Spafford, J. S. Meredith, and J. S. Vetter, “Quantifying NUMA and contention effects in multi-GPU systems,” in *Proceedings of the Fourth Workshop on General Purpose Processing on Graphics Processing Units*, 2011, pp. 1–7.
- [7] N. R. Tallent, N. A. Gawande, C. Siegel, A. Vishnu, and A. Hoisie, “Evaluating on-node GPU interconnects for deep learning workloads,” in *International Workshop on Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems*. Springer, 2018, pp. 3–21.
- [8] A. Li, S. L. Song, J. Chen, J. Li, X. Liu, N. R. Tallent, and K. J. Barker, “Evaluating modern GPU interconnect: PCIe, NVLink, NV-SLI, NVSwitch and GPUDirect,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 1, pp. 94–110, 2019.
- [9] C. Pearson, A. Dakkak, S. Hashash, C. Li, I.-H. Chung, J. Xiong, and W.-M. Hwu, “Evaluating characteristics of CUDA communication primitives on high-bandwidth interconnects,” in *Proceedings of the 2019 ACM/SPEC International Conference on Performance Engineering*, 2019, pp. 209–218.
- [10] T. Allen and R. Ge, “In-depth analyses of unified virtual memory system for GPU accelerated computing,” in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 2021, pp. 1–15.
- [11] S. Chien, I. Peng, and S. Markidis, “Performance evaluation of advanced features in CUDA unified memory,” in *2019 IEEE/ACM Workshop on Memory Centric High Performance Computing (MCHPC)*. IEEE, 2019, pp. 50–57.