

MAGNUM IO GPUDIRECT STORAGE (GDS)

Kushal Datta | STORAGE DEVELOPER CONFERENCE, SEP 28-29, 2021

AI #EXTREMESCALE Exploding Data in HPC, Cloud and Enterprise



BIG DATA GROWTH 90% of the world's data in the last 2 years

Source for Big Data Growth chart: IDC - The Digitization of the World (May 2020) storagenewsletter.com: 64.2ZB of Data Created or Replicated in 2020

GROWTH IN SCIENTIFIC DATA Fueled by accurate sensors and simulations



EXTREME-SCALE AI REQUIRES EXTREME IO



LATENCY IO latency dominates as compute shrinks

Extreme throughput required as microservices scale across data center



180 Industry-Leading Organizations



2,500 Active Contributing Members

COMPATIBILITY

Diverse set of standards and solutions in the ecosystem



GDS Architecture and cuFILE

Benchmarks and Use Cases

Ecosystem and CTA

GETTING DATA TO GPUs: IO

WITHOUT GPUDIRECT STORAGE



Low Bandwidth | High Latency | Limited Capacity

WITH GPUDIRECT STORAGE



Higher Bandwidth | Lower Latency O(PB) capacity | CUDA programming model



GPUDIRECT STORAGE(GDS) ADVANTAGE





NFSoRDMA NVMe-oF/RDMA



APPLICATION CUFILE USER API Enduring API for applications and frameworks nvidia-fs driver API For file system and block IO drivers Vendor-proprietary solutions: no patching avoids lack of Linux enabling

NVIDIA is actively working with the community on upstream first to enable Linux to handle GPU VAs for DMA

GPU Memory

Storage DMA Engine

OS KERNEL or THIRD-PARTY ALTERNATIVE



GPUDIRECT STORAGE: DYNAMIC ROUTING



GDS, affinity	$\underline{NVMe}_A \rightarrow GPU_A$ Target
Non-GDS, affinity	$\underline{NVMe}_{A} \to CPU_{A} \operatorname{\underline{Buf}} \to GPU_{A} \operatorname{\underline{Target}}$
NVLink, no affinity	$\underline{\text{NVMe}}_{\text{A}} \rightarrow \text{GPU}_{\text{A}} \operatorname{Buf} \rightarrow \text{GPU}_{\text{B}} \operatorname{Target}$
Non-GDS, no affinity	$NVMe_A \rightarrow CPU_A Buf \rightarrow GPU_B Target$
GDS, no affinity	$\underline{NVMe}_{A} \to \mathbf{GPU}_{A} \operatorname{\underline{Buf}} \to \mathbf{CPU} \to \mathbf{GPU}_{B}$

4 Gen4 drives DGX A100

Best overall

- Slower: bounce buffer
- Best for cross-socket
- Agnostic to affinity thru CPU
- Some overheads



COMPATIBILITY MODE No Appreciable Degradation Relative to Non-GDS at 128+KB

- cuFile APIs remain functional when:
 - GDStorage-enabled drivers and nvidia-fs.ko are not installed
 - Mounted file system is not GPUDirect Storage compatible
 - File system's specific conditions for O_DIRECT are not met
 - Development is without root privileges
- Configurable in config.json by admin or user
- Advantages:
 - Enable code development on platforms without GPUDirect Storage support
 - Containerized deployment across a range of systems
- cuFile vs. POSIX APIs for moving data to GPU memory:
 - Performance is on par: slight overhead at smallest IO sizes
 - Measured results on DGX-2 and DGX A100 with local storage







CUFILE API

cuFILE LIBRARY

Features

- Shared user-space file IO library for GDS-ready DMA and RDMA file systems
- Supports aligned and unaligned file IO offsets and sizes
- Avoids cudaMemcpy, intermediate copies to page cache, and user-allocated sys memory
- User-space RDMA connection and buffer registration management
- Dynamically routes IO using optimal path based on HW topology and GPU resources
- Supports compatibility mode that works with user-level installation only for development

Requirements

Needs files to be opened in O_DIRECT mode to bypass page cache and buffering in system memory for POSIX-based file systems

Documentation

https://docs.nvidia.com/gpudirect-storage/overview-guide/index.html



USAGE EXAMPLE Write from GPU Memory to Local or Remote Storage

```
int main(void) {
                                                               // set GPU device, etc.
...
CUfileError t status;
CUFileDescr t cf descr;
                                                               // general enough to support Windows
CUFileHandle t cf handle;
status = cuFileDriverOpen();
                                                               // initialize
fd = open(TESTFILE, O WRONLY|O DIRECT, 0644);
                                                               // interop with normal file IO
cf descr.handle.fd = fd;
status = cuFileHandleRegister(&cf_handle, &cf_descr);
                                                               // Check support for file at this mount
cuda result = cuMemAlloc(&devPtr, size);
                                                               // User allocates memory
status = cuFileBufRegister(devPtr, size);
                                                               // performance optimization
assert(cudaMemset((void *) devPtr, 0xab, size) == cudaSuccess);
ret = cuFileWrite(cf handle, devPtr, size, 0, 0);
                                                               // ~pwrite: file handle, GPU base address, size,
                                                              // offset, GPU buffer offset
status = cuFileBufDeregister(devPtr);
                                                               // optional cleanup for good hygiene
cuFree(devPtr);
cuFileHandleDeregiser(cf handle);
close(fd);
cuFileDriverClose();
return 0;
```



HOW TO DEVELOP WITH GDS? What SW Is Needed to Get GDS to Work with DGX A100



***ONE CHANGE FITS ALL** EXT4 with local (NVMe) or direct attached storage (NVMe-oF) NFS (VAST, Isilon/PowerScale) | DFS (EXAScaler, WekaFS)

MECHANISM OF ACTION

Install GDS (UML + KMD) Code change required to use cuFile API*

Install GDS (UML + KMD) Import DALI reader (open-source)*

> Install GDS (UML + KMD) **GDS** integrated





BENCHMARKING GDS

GPUDIRECT STORAGE - BENEFITS Higher Bandwidth, Lower Latency, Lesser CPU utilization



Higher bandwidth - Direct path leads to more throughput

Low latency with GPUDirect Storage Fairly flat (Predictable) Peak IO bandwidth using fewer CPU cores for IO



DGX A100 LOCAL STORAGE COMPARISON GDS advantage is capped by limited throughput of local drives

2.8x performance advantage at 4K for no copy shrinks to 1x when amortized on larger IO sizes

3-5x CPU utilization advantage

4 drives @ 26 GiB/s saturate 1 GPU



IO Size



NVIDIA DGX A100 SYSTEM DIAGRAM





GDS CONFIGURATIONS ON DGX Lower Latencies and Increased CPU Efficiency

STANDARD GDS CONFIGURATION IN DGX A100



North-south storage configuration adapters 4 and 5 used for GDS



2 N-S NICS - GDS VS CPU-GPU (READ)

- Use both E-W NIC on slot 4 (socket 1) and slot 5 (socket 0). •
- Use GPU 7 (socket 1) and GPU 3 (socket 0) for best perf. •
- GDS has up to 1.47x in BW advantage over CPU-GPU. •
- GDS has up to 3x in CPU utilization advantage.
- BW achieved with 2 E-W NICs is nearly 2x of 1 E-W NIC. •



GDS Advantage



GDS CONFIGURATIONS ON DGX Lower Latencies and Increased CPU Efficiency



North-south storage configuration adapters 4 and 5 used for GDS

Converged east-west configuration adapters 0-3, 6-9 used for GDS



IO BENCHMARK ON DGX A100

WITH GPUDIRECT STORAGE



Performance benchmarking done with the gdsio tool using standard GDS configuration in DGX A100 on N-S network adapters, UB 20.04 MLNX_OFED 5.3 GDS 1.0

NVIDIA DGX A100 SYSTEM





GDS ENABLES PEAK IO BANDWITH ON DGX A100

WITHOUT GPUDIRECT STORAGE





Performance benchmarking done with the gdsio tool using standard GDS configuration in DGX A100 on N-S network adapters, UB 20.04 MLNX_OFED 5.3 GDS 1.0

WITH GPUDIRECT STORAGE





CASE STUDIES

CASE STUDY 1: ACCELERATING DL INFERENCE WITH GDS GDS-Specific Gains Are up to 2.6X with DALI, 1.9X Without



Performance benchmarking done with the gdsio tool using standard GDS configuration in DGX A100, UB 20.04 MLNX_OFED 5.3 GDS 1.0. Application for batch size >= 32 limited by GPU compute throughput



CASE STUDY 2: ACCELERATING SEISMIC WITH GDS

FORWARD MODELING



BACKWARD MODELING

Compression Quality



LEARNINGS/BEST PRACTICES

- **Relaxed ordering** on the NICs must be to enabled.
- When running gdsio, the proper NUMA node(s) for the GPUs being used need to specified in order to get best performance.
- When using 2 E-W NICs, bind the GPU to NICs' NUMA node respectively for best performance. This way, both NICs can be utilized.
- Need ~96 threads per GPU to max out a 200 Gbps NIC but the latency might be large.
 - **Disable GDS IO path stats** for best performance.
 - echo 0 | sudo tee
 - /sys/module/nvidia_fs/parameters/rw_stats_enabled
- Both client and server need to have the same "peer_credits", "peer_credits_hiw", "concurrent_sends" settings. Both need to increase "lnet_transaction_timeout" to 100. See JIRA.
- Contain accesses to the same socket or even the same PCIe tree. Avoid writes across CPU Root Ports.





✓ LIQID



GPUDirect FAMILY WITHIN MAGNUM IO A natural progression

GPUDirect peer to peer - enables access between GPU peers on PCIe, NVLink Used by MPI, UCL, NCCL, NVSHMEM

GPUDirect RDMA - enables direct DMA data path from NIC to GPU memory - no CPU Used by MPI, UCL, NCCL, NVSHMEM, GPUDirect Storage

GPUDirect Storage - enables direct DMA data path from storage to GPU memory Used by applications, frameworks, and readers



MAGNUM IO TECHNOLOGIES Historical and New Offerings in Data Movement, Access, and Management

NVIDIA Innovations | Peak Performance | RDMA Everywhere | Flexible Abstraction





CALL TO ACTION

EXPLORE: https://developer.nvidia.com/gpudirect-storage

GPUDIRECT DEVELOPER BLOG: https://developer.nvidia.com/blog/gpudirect-storage/

> NVIDIA MAGNUM IO ON GITHUB: https://github.com/NVIDIA/MagnumIO





GDS Engineering

GDS Product

Maitree Kanungo

GDS Architecture

CJ Newburn, Kiran Kumar Modukuri

Ecosystem

Barton Fiske, Michal Shoham, Rob Davis

Management

Ira Nanda

Sandeep Joshi, Aniket Borkar, Rebanta Mitra, Sourab Gupta, Zhen Zeng, Vahid Noormofidi



