

# IBM announces new performance for ESS 3200 with NVIDIA DGX & GPUDirect Storage

Silverton Consulting, Inc. StorInt™ Briefing

IBM® recently announced new performance results with IBM Elastic Storage® System (ESS) 3200 on NVIDIA DGX™ systems using the GPUDirect™ Storage interface.

## IBM ESS 3200 DGX A100 GPUDirect performance

The ESS 3200 is IBM's latest AFA storage appliance running IBM Spectrum® Scale storage software to provide high performance file system storage that supports up to 80GB/sec per node, which is linearly scalable in performance for up to 10TB/sec and in capacity for up to Yottabytes (YB, millions of PB) of storage.

NVIDIA DGX systems are a new AI data center infrastructure, that is designed for ease of use and deployment of AI technologies in enterprise and HPC environments. NVIDIA DGX A100 systems are configured in PODs, with 1 to 4 DGX A100 nodes per rack. NVIDIA also offers the DGX SuperPOD, that can scale from 20 to 140 DGX A100 nodes to support almost supercomputing levels of performance.

A DGX A100 system has 8 NVIDIA A100 (40 or 80 GB) GPUs, dual AMD 64 core CPUs, with 1 or 2TB of memory, internal NVMe SSD storage, 8 external compute and 2 external storage HDR (200Gbs) InfiniBand interfaces for compute and storage fabrics.

NVIDIA GPUDirect Storage is a new storage interface designed to speed up data transfers to GPU memory. Prior to GPUDirect, reading data from storage and placing it in GPU memory would require the data to be moved to CPU memory before it could be sent to a GPU. With GPUDirect Storage, data can be read directly from internal or external storage to GPU memory without having to stop in CPU memory. Writes from GPU memory to storage follow a similar path.

IBM had benchmarked their ESS 3200 storage for DGX A100 (not using GPUDirect) and found it could sustain decent read throughput. But when they redid their benchmarks using the GA version of NVIDIA's Magnum IO stack with GPUDirect Storage with a single node DGX A100 POD, a single ESS 3200 system was able to deliver **43GB/s read bandwidth to the DGX A100's 8 GPUs or 1.9X better throughput than non-GPUDirect Storage, over the storage fabric.** This level of performance is ~86% of the theoretical read bandwidth available for the InfiniBand storage fabric.

To further stress ESS 3200 storage performance, IBM created a special IO benchmark that used the (8 HDR InfiniBand) DGX compute fabric (not a supported NVIDIA configuration) to perform read IO. In this case, they configured a pair of ESS 3200 storage systems and two DGX A100 nodes to run the FIO read benchmark, modified to use GPUDirect Storage (Beta version) and the DGX compute fabric. Here, the pair of ESS 3200s were able to provide **191.3GB/s of read bandwidth to the pair of DGX A100 nodes 16 GPUs across the 8 HDR InfiniBand compute fabric.** The pair of ESS 3200s effectively saturated the DGX A100 compute InfiniBand fabric. IBM and others have used this approach to measure max throughput for DGX A100 GPUs even though it's not a supported NVIDIA configuration.

IBM has also updated their reference architecture<sup>1</sup> describing how best to deploy ESS 3200 storage with 1, 2, 4 or 8 A100 node DGX PODs.

## IBM ESS 3200 and NVIDIA DGX SuperPOD certification

IBM also announced that they would be working with NVIDIA to certify ESS 3200 storage for use with NVIDIA DGX SuperPOD. Given the inherent scalability of the ESS 3200 system, it would seem the ideal storage for NVIDIA DGX SuperPOD deployments.

As discussed earlier, NVIDIA DGX SuperPOD comes in units of 20 DGX A100 nodes, with all the networking and other hardware to support AI activity at this level and can scale up to 140 DGX A100 nodes.

NVIDIA envisions that a DGX SuperPOD would provide a core AI capability for the enterprise/lab and DGX PODs could supply edge AI infrastructure for inferencing and real-time (re-)training.

## Significance

AI seems to be becoming a central component in more new enterprise and HPC workloads. However, the infrastructure required to sustain AI capabilities for these workloads is specialized and expensive to use in isolation. And scaling this infrastructure up to support more AI activity can be challenging.

NVIDIA, with their DGX POD and SuperPOD systems, based on DGX A100 nodes, InfiniBand networking and scalable storage have made supplying shared, data center class, AI infrastructure for enterprises and HPC labs much simpler to purchase, deploy and implement.

However, often the weak link in AI activity is data. AI model training can involve multiple passes over vast datasets, to attain needed accuracy. Storage that can support this AI data IO can be difficult to find. And storage that can scale to the size and performance needed to keep DGX POD or SuperPOD GPUs well supplied with data, has been nigh impossible.

But IBM ESS 3200 AFA storage and NVIDIA Magnum IO GPUDirect Storage has been shown the capability to support this work, at the level needed to keep 10s to 1000s GPUs busy and DGX storage fabrics humming. Once ESS 3200 DGX SuperPOD certification is achieve, it will be relatively straight forward to deploy ESS storage with DGX SuperPOD.

In the meantime, for smaller DGX POD environments, the IBM Reference Architecture can be used as a recipe as to how to configure ESS 3200 storage for use with 1, 2, 4 or 8 DGX A100 nodes.

---

***Silverton Consulting, Inc., is a U.S.-based Storage, Strategy & Systems consulting firm offering products and services to the data storage community.***

---

<sup>1</sup> Please see <https://www.ibm.com/downloads/cas/MJLMALGL> as of 26 June 2021