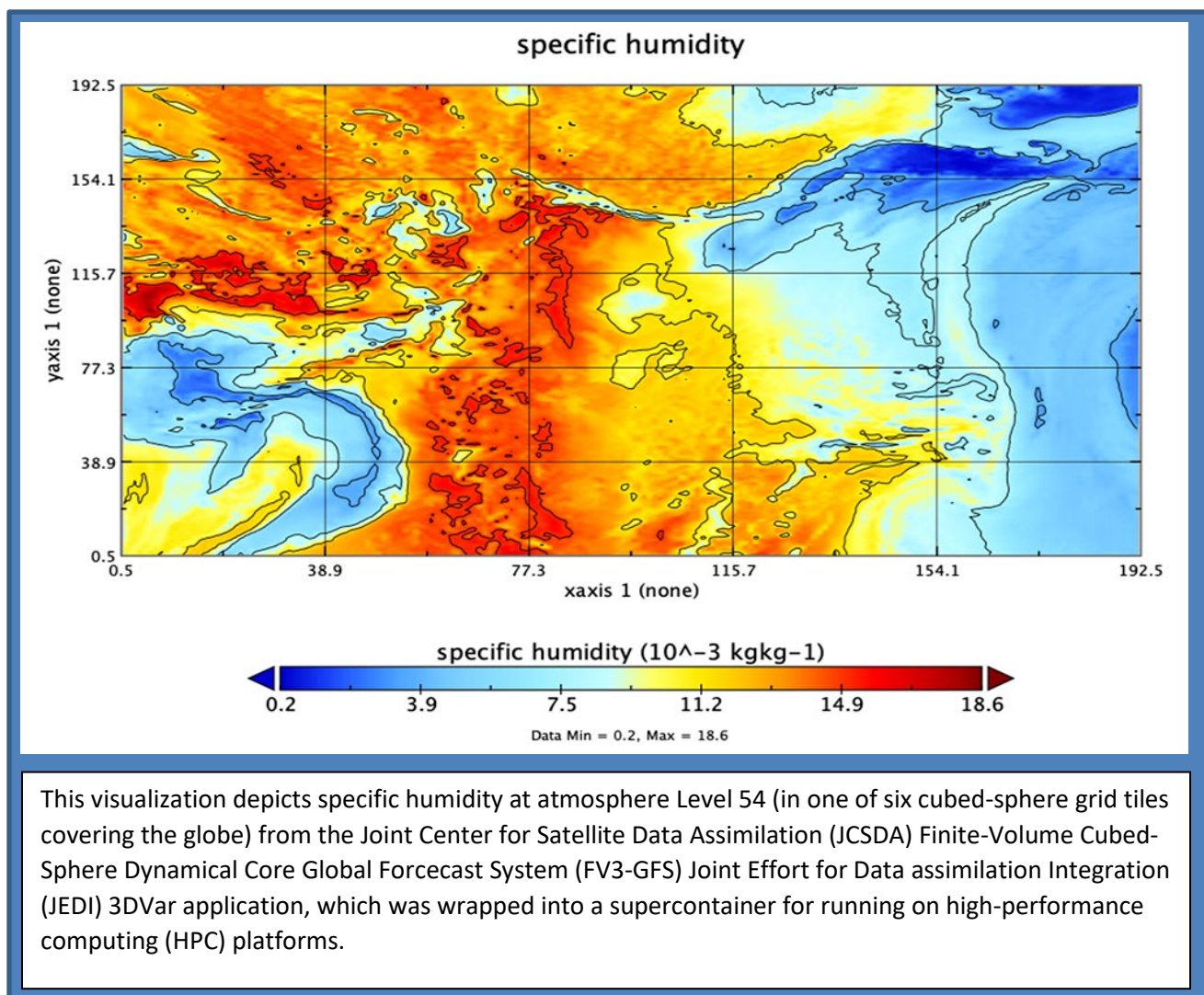


# Building Storage Solutions for High-Performance Computing

## Kioxia, NVIDIA, Weka



High-Performance Computing (HPC) originally started out in defense and research domains, such as universities and government agencies. The workloads that the initial HPC solutions were used for included particle physics, fluid dynamics, nuclear weapons modeling, and research and aerospace projects – and, obviously, the space program. What we have seen in the last several years is HPC being applied to new problems in commercial enterprises. This includes product modeling to speed time to market, business intelligence (analyzing what customers are doing), and for artificial intelligence and machine learning. These and other areas have formed a toehold for HPC in the commercial space and enterprise, which will undoubtedly grow over time.



As you can imagine, there are a lot of differences between enterprise datacenter storage architectures and HPC storage architectures. Conventional datacenters (especially private cloud implementations) are really built to simplify the job of the IT staff by using homogeneous resources that simplify migration and scaling and flexibility. Virtualization of all resources (compute, storage, and networking) is really important to achieving application availability and data integrity, with a focus on avoiding loss of data and downtime. Performance, generally, while it can be important, is secondary to flexibility and data integrity concerns. While there are some enterprise architectures implemented for performance such as data lakes, these are outliers.

The priorities for HPC architectures are the opposite of enterprise datacenters architectures. You might have thousands or hundreds of thousands of processors and GPGPUs that are drawing data from a single storage pool. They tend to run applications as projects or batches, where a job is set up, ran, and torn down when complete. Because of that, HPC architectures are highly performant and they don't necessarily have to be standard. Being able to configure a cluster for optimum performance is usually the primary concern.

The challenges of moving HPC into the enterprise require IT staff to think very differently. You have to plan differently, manage resources differently, and, in the end, you have to execute very differently because this is a very different environment. Maximizing HPC storage performance was the subject of the G2M Research February 23<sup>rd</sup> webinar, with Kioxia, NVIDIA, and Weka.

[Kioxia](#) discussed the importance of local storage in each compute node to maximize performance by reducing storage latency, especially for extremely large datasets, and the importance of connecting storage nodes to compute nodes via NVMe-oF™. Kioxia also talked about the fit of their various SSDs for HPC architectures.

[NVIDIA](#) spoke about using InfiniBand as the interconnect within supercomputing clusters, as it is the most widely deployed interconnect in HPC due to its deterministic performance. They also discussed using NVIDIA's GPUDirect Storage (MagnumIO) interface to avoid moving data through the CPU.

[Weka](#) talked about I/O patterns and how they vary across HPC workloads, how GPUs have increased the density of HPC clusters, and why current NAS solutions cannot provide the bandwidth required to feed these clusters.

Kioxia, NVIDIA, and Weka provide approaches to storage architectures that maximize the performance of HPC Clusters. Hear the discussion from our last webinar [here](#).



Mike Heumann, Managing Partner  
G2M Communications