

CCC BUILDS A SUPERIOR AI INFRASTRUCTURE WITH EXCELERO'S NVMeSH STORAGE SOFTWARE

Eliminates GPU storage bottleneck, to achieve 3-4x faster analysis of ML training datasets

CASE STUDY



INTRODUCTION

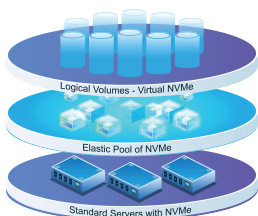
Founded in 1980, CCC Information Services (CCC) is an innovative technology provider to the automotive, insurance, and collision repair industries. Its powerful CCC ONE™ platform connects a vast network of 350+ insurance companies, 25,000+ repair facilities, OEMs, thousands of parts suppliers, and dozens of third-party data and service providers. CCC's smart, flexible, and intuitive solutions make automotive data actionable - enabling customers to make informed decisions that deliver faster and better experiences for end consumers.

For example, CCC® Smart Estimate (Smart Estimate) is the world's first in-production AI solution for vehicle collision estimates. It leverages CCC's estimating logic and AI photo analytics capabilities to pre-populate an estimate with suggestions for human estimators to review, edit, and advance. The solution, live in market with insurance customers, allows users to upload mobile phone images of collision damage along with other attachments – then applies CCC's estimating logic and AI to vehicle collision photos to predict repair requirements and suggest estimate lines, including parts likely required to complete the repair.

Training the machine learning (ML) model in the Smart Estimate workflow, which can include a series of AI-based solutions, required a computing infrastructure with massive scalability. It was a given that the infrastructure needed GPU servers and high-speed networking infrastructure. When it came to storage, however, CCC needed a scale-out Flash storage option in order to deliver the high throughput and low latency the workflow required. By deploying Excelero's NVMeSH®, a software-defined all-NVMe scale-out block storage solution, to supercharge the BeeGFS distributed parallel-file system from ThinkParQ, CCC achieved 3-4X times faster processing of its ML datasets. The time-savings allowed CCC's team of AI developers to train more datasets, faster – while assuring exceptional performance for its customers.

Benefits of Elastic NVMe for AI and Machine Learning Elastic NVMe for GPUs

Elastic NVMe for GPUs



- Share NVMe resources across multiple GPU servers
- Access remote NVMe at local speed
- Exceed the performance and capacity limits of local flash on GPU servers
- Eliminate the need to copy data locally conserving time and drive endurance
- Datasets can be larger than what can fit inside the GPU Server
- Zero-CPU storage-target with Excelero's patented Remote Direct Drive Access

CCC BUILDS A SUPERIOR AI INFRASTRUCTURE WITH EXCELERO'S NVMeSH STORAGE SOFTWARE

Eliminates GPU storage bottleneck, to achieve 3-4x faster analysis of ML training datasets

FEEDING THE GPU SERVERS, NOT STARVING THEM

Through the use of machine learning techniques, Smart Estimate AI and human estimators create a continuous self-improving process, allowing fast and ever-smarter auto physical damage (APD) estimates over time. Making this a reality required an innovative approach to an AI-centric data center architecture. CCC built its next revision of GPU data center cluster based on NVIDIA DGX GPU servers across a high-speed 100 Gigabit Ethernet (100GbE) networking backbone. The companion storage systems and software however required a rethinking of the typical approach.

To date, CCC's AI logic behind the application has ingested millions of small file data sets. Often data ingestion took place at a pace of millions files every few hours. Such a massive volume of small files is a classic challenge for any storage solution, because it demands high throughput and low latency – often the opposite of what traditional storage systems provide.



MASSIVE VOLUME OF SMALL FILES: A STORAGE CHALLENGE

Such intensive processing can also create a “GPU bottleneck” often found in AI deployments – where GPU systems can’t ingest the data fast enough, and user response times slow. With the help of system integrator Advanced HPC, CCC evaluated scale-out storage designs.

CCC knew that deploying a parallel distributed file system was key to handling the metadata of a large number of small files efficiently, avoiding bottlenecks. It knew its existing controller-based storage resource would need to be bolstered for its throughput and latency requirements, and that newer NVMe SSDs were needed. Lastly, CCC placed NVMe storage on a software-defined platform to enable flexible choices in underlying hardware and make it available across any network as if it were local storage – to avoid hardware vendor lock in, boost efficiency and ROI.

CCC trialed the Excelero NVMeSH block storage software which is purpose-built for AI and ML deployments. Because NVMeSH ensures efficient use of both the GPUs themselves and the associated NVMe flash in CCC's Smart Estimate deployment. The end result is higher ROI, easier workflow management and faster time to results. Excelero and Advanced HPC assisted in optimizing system-wide tuning to assure optimal network performance across the existing network fabric.

In implementing Excelero NVMeSH® software-defined block storage along with the BeeGFS parallel distributed file system, CCC achieved 3-4X times faster analysis of the ML datasets. Its system handles 1.5 million IOPs per client node, a measure of storage performance – roughly 10x to 15x its previous capability. The time-savings allowed CCC's team of AI developers to train its model, faster – while keeping exceptional precision, and helping assure the AI-backed collision estimating application delivered value to end consumers in timely manner.

“Using Excelero NVMeSH along with the BeeGFS parallel file system gave us three to four times faster analysis of ML data sets. Throughput is now 10 to 15 times greater, and latency is negligible.”

Andrey Ptashnik, Lead Enterprise Architect, CCC

CCC BUILDS A SUPERIOR AI INFRASTRUCTURE WITH EXCELERO'S NVMesh STORAGE SOFTWARE*Eliminates GPU storage bottleneck, to achieve 3-4x faster analysis of ML training datasets*

"To advance CCC's vision of 'AI Everywhere' we continue to invest and deploy tools that support our world-class AI/ML solution development infrastructure. The Excelero technology and NVMesh product are important elements in our broader environment."

Reza Rooholamini, VP, Chief Architect, CCC

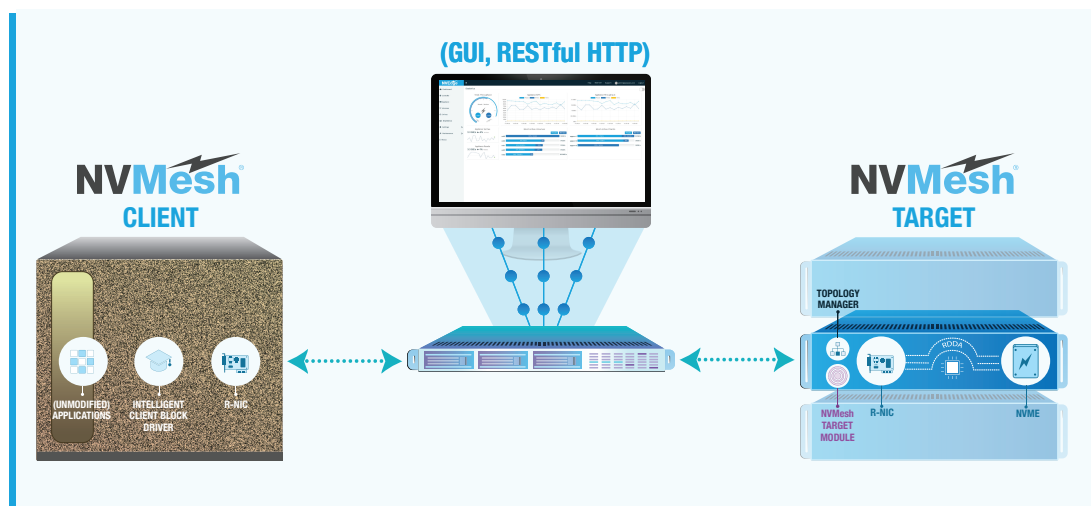
OPTIMIZING AI IN THE DATA CENTER

Smarter approaches to the GPU storage bottleneck have evolved today, helping resource-hungry AI and ML applications achieve the high-throughput and low latency required for superior application performance to end users.

The biggest advantage of modern GPU computing is also creating its biggest challenge: GPUs have an amazing appetite for data. Current GPUs can process up to 16 GBs of data per second. The DGX-1 has a theoretical limit of 7.8 GB/s bandwidth, but with only 4 SATA SSDs, it is limited to about 2.2 GB/s.

Theoretically, it can process 2 million random IO/s, but local storage only provides 400K IO/s. NVIDIA's latest DGX-2™ system has as many as 16 GPUs, but with 30 TBs (8 x 3.84TB) of

in-server NVMe drives, this is by far not enough local storage. Other brand GPU servers typically feature few PCIe lanes for NVMe drives, meaning that the lowest latency option for these servers is a severe bottleneck or simply has too little capacity for the required datasets. Starving the GPUs with slow storage or wasting time copying data wastes expensive GPU resources and affects the ROI.



"At Advanced HPC we're always searching for innovative solutions that save our customers time and money. Excelero NVMesh is a total game changer for AI - one of those products that I've truly only seen once or twice over my 20-year storage career. Its small block I/O is particularly impressive, as CCC experienced, and it delivers unprecedented IOPS and throughput at near-zero latency to GPU servers. Our customers are doing 5-10x the work by adding just one of our 2U appliances powered by Excelero's NVMesh and our world-class support and engineering staff."

Joe Lipman, Senior Sales Engineer, Advanced HPC, who led the CCC deployment

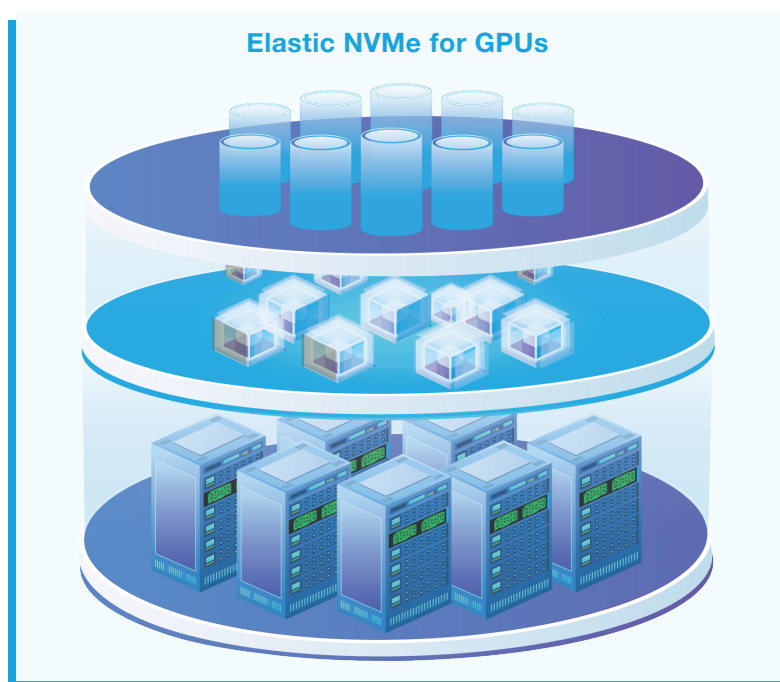
CCC BUILDS A SUPERIOR AI INFRASTRUCTURE WITH EXCELERO'S NVMesh STORAGE SOFTWARE

Eliminates GPU storage bottleneck, to achieve 3-4x faster analysis of ML training datasets

SUMMARY

The incredible capability of GPUs and the rise of affordable compute power, challenges IT teams to think at data center scale, thus leveraging the ability to apply AI, ML and Deep Learning techniques to large data pools, while making sure the entire system is scalable, highly performant and efficient.

GPUs have become the go-to compute resources behind AI and ML workloads and NVMe flash has become the standard for high-performance, low latency storage. But AI and ML workloads need much more capacity than that available locally in GPU systems. NVMesh enables users to share NVMe resources across GPU servers, such as Nvidia DGXs, or add unlimited external NVMe storage with local performance.



Excelero NVMesh enables AI data scientists and HPC researchers to feed huge amounts of data to their applications. They no longer face data bottlenecks and can get to better results faster.

NVMesh is an Elastic NVMe storage software solution that provides GPU systems with access to massively scalable, extreme-performance NVMe flash storage as if it were local flash. The end result is higher ROI for both GPUs and NVMe, easier workflow management and faster time to results.

NVMesh Features for GPU

- **NVMesh unifies remote NVMe devices into a logical block pool that performs the same as local NVMe flash**
- **NVMesh allows full utilization of the IO/s and bandwidth capabilities of NVMe drives across a network**
- **Nvidia DGX-1 and DGX-2 can use their massive network connectivity to access remote NVMe logical volumes, with redundancy if desired!**
- **MUCH faster than local SATA SSDs**
- **Larger shared pools than possible within the platform**
- **Other GPU optimized systems can access remote NVMe at local latencies and bandwidth**
- **Random IO characteristics of NVMe preserved, achieving 10s of millions of potential IO/s at minimal latencies**