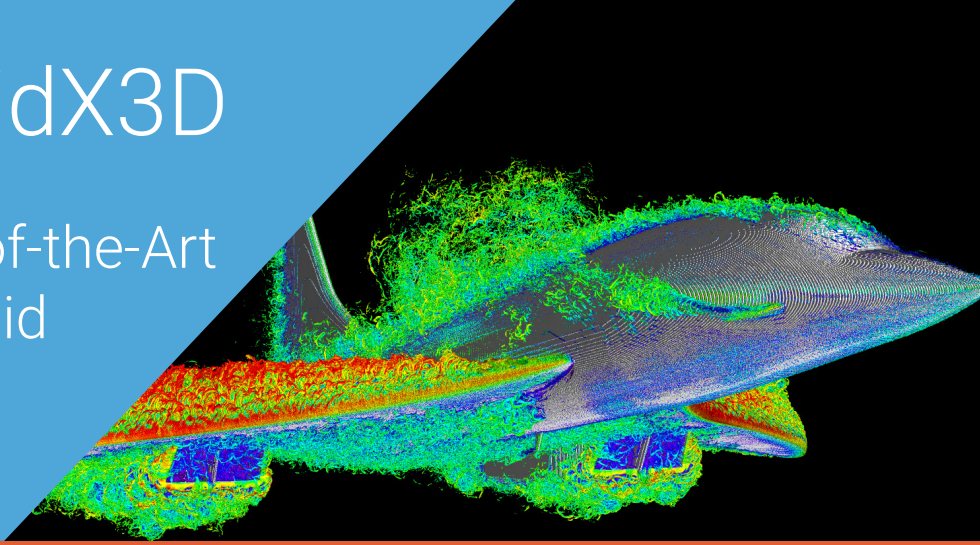




FluidX3D

Advancing the State-of-the-Art for Computational Fluid Dynamics (CFD)

CASE STUDY



Innovative Approaches to Solving for Extreme-Scale CFD

What You Get with SuperNODE for CFD

BREATHKING PERFORMANCE

Up to 32 accelerators on a single server with terabytes of VRAM to solve massive problems

ULTIMATE FLEXIBILITY

Easy reconfiguration of heterogeneous accelerators from one to several

OUTSTANDING SIMPLICITY AND EASY MANAGEMENT

Turnkey, easy to deploy, no software changes needed

COMPELLING ECONOMICS

Lower TCO, higher resource utilization, and less complexity, space, power, and cooling

END-TO-END SOLUTION

Integrated computing, rendering, and storage

CFD has revolutionized the design and development process in the aerospace, automotive, process and chemical, and other industries. It can significantly reduce product development costs and risks, improve functional performance, and model many physical phenomena like weather, groundwater flows, and plasma physics.

The Navier-Stokes equations describe fluid dynamics very well. However, accurately solving these equations at scale remains daunting, limited by the computational capabilities and costs of resolving the smallest space and time features.

Dr. Moritz Lehmann, the sole developer of FluidX3D, a fast and highly memory-efficient CFD software based on the Lattice Boltzmann Method (LBM), recently solved one of the largest CFD problems ever. This simulation of one second of a Concorde landing at 300 km/h with a resolution of 40 billion cells (*Figure 1*) took 33 hours on a GigalO SuperNODE™ system with 32 AMD Instinct™ MI210 accelerators. Conventional CFD programs could take years to solve this problem.

By running FluidX3D on the GigalO SuperNODE system, which supports up to 32 accelerators or Graphics Processing Units (GPUs) and provides terabytes of directly accessible fast Video Random-Access Memory (VRAM), engineers can rapidly solve these massive CFD problems at scale. Moreover, they can quickly visualize the results because FluidX3D also allows for rendering massive simulation data directly on the same VRAM as the GPU accelerators supported on the GigalO SuperNODE.

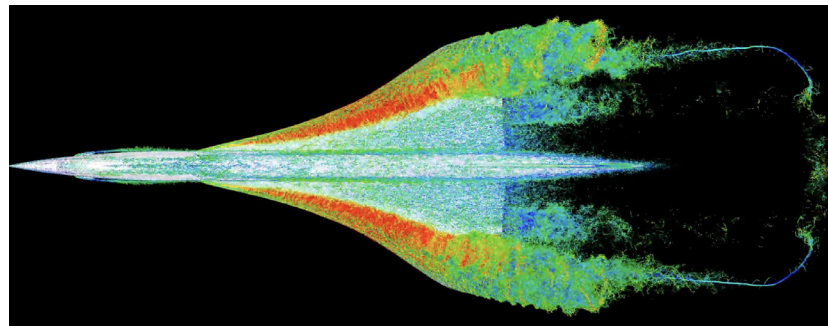


Figure 1:
Concorde Landing Simulation

“Over the weekend, I tested FluidX3D on the world's largest HPC GPU server, GigalO SuperNODE. I ran one of the largest CFD simulations ever, the Concorde for 1s at 300km/h landing speed. 40 Billion cells resolution. It took 33 hours to run on 32 GPUs.”

— Dr. Moritz Lehmann, Sole Developer of FluidX3D

Why Is the GigaIO SuperNODE Exceptional for Extreme-Scale CFD?

Novel algorithms and computing architectures are required to solve extreme-scale CFD problems. The GigaIO SuperNODE is ideal for these workloads because it is the only high-performance system that:

- Supports up to 32 GPUs on a single server, so CFD applications that perform exceptionally well on GPUs and run parallel across these GPUs can scale to exceptional levels.
- Delivers unmatched capability to directly address and leverage all available VRAM (several terabytes) so enormous problems can fit in memory, making the CFD application run extraordinarily fast, as communication happens through fast PCIe connections instead of Ethernet networking.
- Offers a single, dynamic, open, industry-standard (PCIe) high-performance network fabric (FabreX™) for speedy and seamless access to all computing and storage resources (other systems typically contain two or more network fabrics, adding unnecessary complexity). FabreX minimizes communication overhead during parallel computing and allows the fast end-to-end coupling of computing, rendering, and storage on one GigaIO SuperNODE (Figure 2).

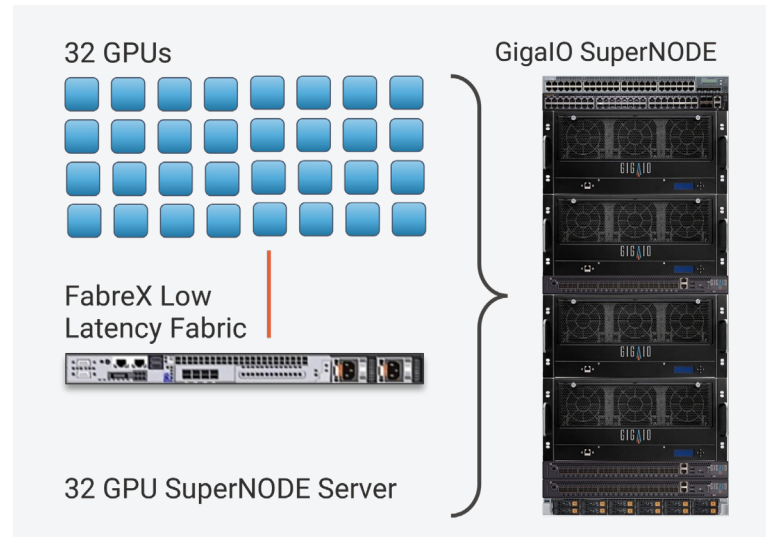


Figure 2: The GigaIO SuperNODE

This flexibility is a game-changer for CFD. The SuperNODE can support up to 32 high-performance GPUs in a single system, providing an unprecedented amount of parallel processing power and video memory to tackle the largest, most complex fluid dynamics simulations. And because the GPUs are dynamically connected via low-latency, high-bandwidth FabreX fabric, data can be seamlessly passed between them with minimal overhead, enabling exceptional scaling.

The advantages of the SuperNODE extend well beyond just raw GPU horsepower. Its composable architecture allows CFD specialists to easily integrate and optimize complementary resources such as fast storage, high-speed networking, and specialized Data Processing Units (DPUs), creating a truly end-to-end solution tailored to their unique requirements.

SuperNODE also boasts the ability to transcend traditional server boundaries, enabling server-to-server communication across the FabreX fabric. This capability opens new frontiers in cluster-scale CFD, allowing individual compute nodes to directly access the system memory of their peers. The implications for massively parallel, extreme-scale simulations are profound.

Many extreme-scale CFD problems are intractable or take too long to generate results. Running FluidX3D on the GigaIO SuperNODE can help clients address many large-scale CFD challenges.

“We are excited to partner with leading innovators like Dr. Moritz Lehmann, the sole developer of FluidX3D, to solve extreme-scale CFD problems. Together, we are advancing CFD’s state-of-the-art to benefit all manufacturing and research organizations.”

— Alan Benjamin, CEO of GigaIO



Overcoming Large-Scale CFD Challenges

Figure 3 lists several challenges in large-scale CFD, along with key elements of solutions to overcome these hurdles.

Some large-scale CFD and IT challenges include the ability to:

- **Handle Complex Flow Scenarios and Intricate Geometries:** Most practical industrial fluid flow problems occur in regions with complex geometries and transient movement, which are hard to mesh in order to resolve the detailed physics with traditional methods.

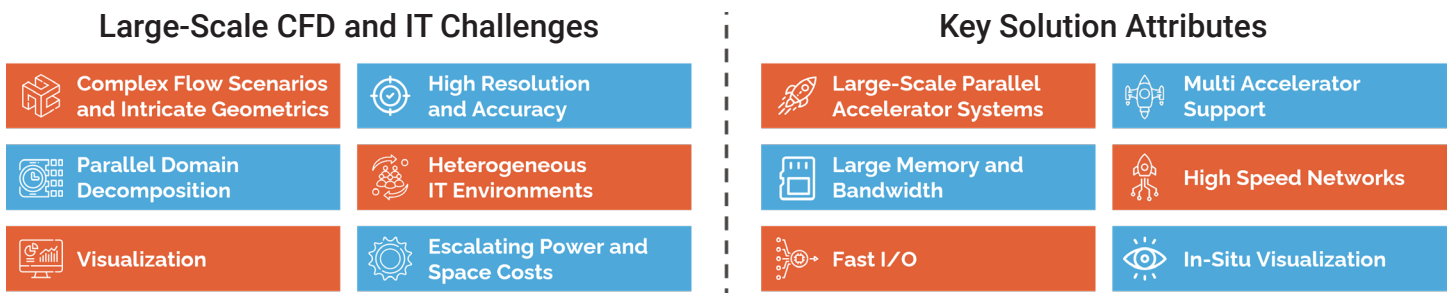


Figure 3: Large-Scale CFD and IT Challenges and Key Solution Attributes

- **Achieve High Resolution and Accuracy:** Traditional CFD algorithms rely on empirical approaches to model turbulence, which may limit resolution and accuracy.
- **Construct Parallel Domain Decomposition:** Most CFD algorithms use domain decomposition techniques to accelerate time-to-results with parallel computing. This is hard to do on complex geometries with moving parts using traditional methods.
- **Support Heterogeneous IT Environments:** The IT infrastructure and components to handle large-scale CFD are moving targets with constant innovation. Clients need the flexibility to deploy the best infrastructure technologies and components to solve their most complex problems.
- **Visualize Simulation Results:** The simulations are typically so large that storing the volumetric data for later rendering becomes unmanageable. If possible, volumetric files should not be exported to the hard disk, as this increases runtime disproportionately.
- **Reign In Escalating Power and Space Costs:** Large-scale CFD typically requires very dense high-performance systems that consume lots of electricity and need special facilities, driving up operating costs.

The key attributes of a system to address the above challenges are:

- **Large-scale Parallel Accelerator Systems:** Many practical CFD problems are growing in size and complexity and need systems containing many high-performance accelerators or GPUs.
- **Multi-accelerator Support:** CFD clients typically use systems that become heterogeneous over time because they regularly upgrade parts of their environment to the latest and best combination of hardware and software. The system must support multiple accelerators in order to protect a client's IT investment.
- **Large Memory and Bandwidth:** As CFD problems get larger, each accelerator must have more memory capacity and bandwidth in order to hold each sub-domain's data without impacting performance.



- **High-Speed Network:** Large-scale CFD problems using domain decomposition must pass data between neighboring sub-domains. A low-latency and high-bandwidth interconnect is needed to minimize communication overheads so that the problems can scale better.
- **Fast I/O:** Transient CFD simulations require high-performance processing power and substantial, fast storage to archive multiple terabytes of data generated at a tremendous rate for analysis and post-processing.
- **In-situ Visualization:** CFD problems are often so large that even the fastest storage and network are inadequate to keep pace with the growth of the generated computational results. Consequently, post-processing becomes a bottleneck to end-to-end simulation performance. One approach to solving this performance imbalance is to reduce the amount of output data by implementing in-situ visualization, which constructs the visualization concurrently with the computing.

Running FluidX3D on GigaIO SuperNODE provides all these desired solution attributes.

FluidX3D Efficiently Solves the Largest CFD Problems

FluidX3D is a speedy, memory-efficient CFD software package implemented using OpenCL™ (Open Computing Language) and intended for GPU acceleration. It addresses many of the challenges stated earlier because FluidX3D uses LBM, which represents a paradigm shift in simulations and enables a broad spectrum of applications previously untenable with traditional Navier-Stokes solvers. For instance, LBM works very well for flow problems on complex geometries with intricate structures like porous media flows. Moreover, LBM delivers outstanding performance on GPUs and other accelerators.

In benchmarks that leverage the GigaIO SuperNODE, FluidX3D runs 100-200 times faster than commercial finite volume solvers on similar hardware. It renders over 1 billion cell simulations on a single GPU in just hours. *Figure 4* depicts results on practical applications across a range of Reynolds numbers using FluidX3D.

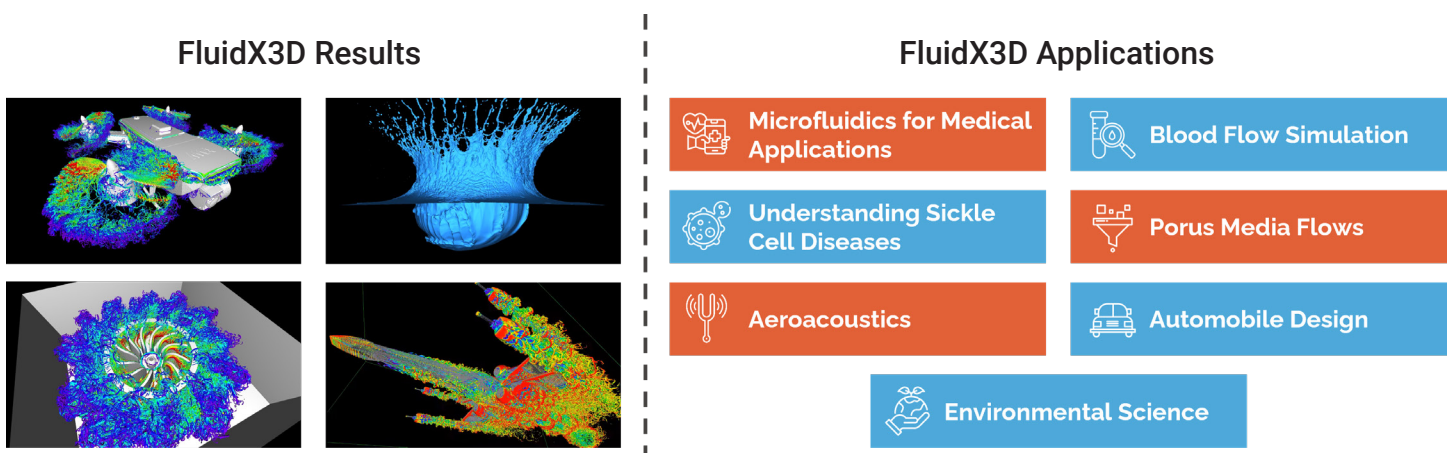


Figure 4: Some FluidX3D Results and Applications

FluidX3D Innovations

FluidX3D also incorporates many innovations to improve the speed of computation while reducing the cost of computation (Figure 5):

- **Memory Optimization:** The program only needs a single grid in memory because it uses a special memory access pattern. This reduces the memory requirement by half.
- **Floating Point Compression:** Using 16-bit floating point memory instead of 32-bit floating point memory cuts the memory requirement by half yet again, without sacrificing the accuracy of the solution.
- **In-situ Visualization:** This allows raw simulation data to be rendered directly in VRAM, so no large volumetric files need to move to the hard disk. The rendering is fully multi-GPU-parallelized via seamless domain decomposition rasterization, and is so fast that it works interactively in real time. With interactive graphics mode disabled, the image resolution can be as large as VRAM allows (4K/8K/16K and above).
- **Supports Many Accelerator Technologies:** Accelerators and GPUs with increased RAM accelerate computational and rendering speed. The OpenCL implementation efficiently permits cross-vendor parallel computing.

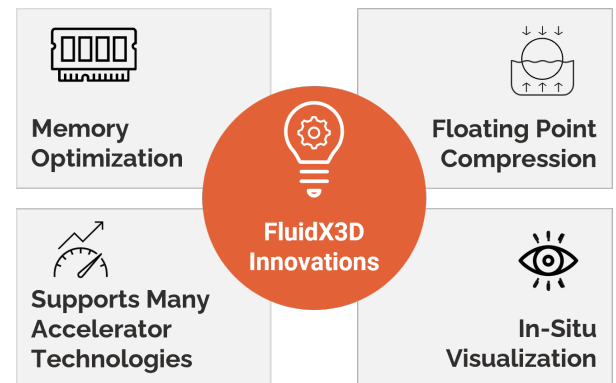


Figure 5: Key Innovations in FluidX3D

These innovations make migrating and optimizing FluidX3D on the GigaIO SuperNODE server easy because its revolutionary new architecture enables a unified, software-defined, and memory-centric composable infrastructure.

GigaIO SuperNODE: The Open Solution for Accelerated Computing

FluidX3D performance is significantly improved when running on the new GigaIO SuperNODE architecture. It maximizes GPU utilization and provides low latency and high bandwidth by integrating compute, GPU acceleration, and fast I/O into a single-system dynamic fabric using standard PCI Express (PCIe) technology.

At the system's heart is FabreX, the dynamic memory fabric that:

- Facilitates connecting memory, network, storage products, and many accelerators to the processors, either directly attached or via server configurations such as NVMe-oF and GPU Direct RDMA (GDR).
- Disaggregates computing, GPUs, storage, and other resource I/O into pools connected by a single-system fabric using PCIe.
- Breaks the server barrier to enable components previously internal to the box to be connected at similar latency and bandwidth but outside the server enclosure.
- Enables server-to-server communication across PCIe and makes cluster-scale computers possible, with direct memory access by an individual server to system memories of all other servers in the network fabric.



These FabreX capabilities drive the following benefits for the GigaIO SuperNODE (Figure 6):

- **Breathtaking Performance:** Delivers the industry's lowest possible latency and highest effective bandwidth. Latency from one server to another is less than 300ns – true PCIe performance across the entire cluster. The FabreX Gen4 implementation scales up to 1 TB/sec bandwidth. In addition, the platform can support up to 32 devices (GPUs, DPUs, FPGAs, etc.), providing room for even higher performance.

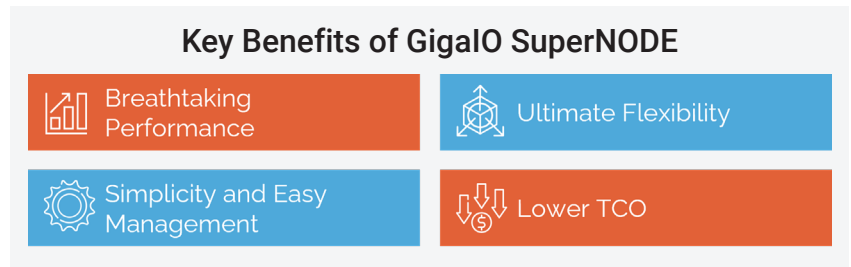


Figure 6: Key Benefits of GigaIO SuperNODE

- **Ultimate Flexibility:** The platform can compose a variety of resources from various manufacturers, such as accelerators (including GPUs, DPUs, TPUs, FPGAs, and SoCs), storage devices such as NVMe, PCIe native storage, and other I/O resources connected to compute nodes. The user is not locked into a single vendor.
- **Outstanding Simplicity and Easy Management:** All GigaIO systems include NVIDIA Base Command Manager in order to provide a complete turnkey solution for supporting workflow-defined infrastructure. This software helps quickly configure and scale compute infrastructure. Even with a heterogeneous mixture of devices manufactured by different vendors, the whole system shows up as a single server with many OpenCL devices.
- **Lower Total Cost of Ownership (TCO):** The system allows clients to upgrade or add compute, GPUs, and application accelerators that plug and play with the user's environment at the component level. Every major subsystem can now operate on its own upgrade cycle. The system's total cost is optimized because FabreX drives a much higher utilization of resources and lowers power, cooling, and space costs.

Conclusions

The GigaIO SuperNODE, in partnership with FluidX3D, represents a significant advancement in solving extreme-scale Computational Fluid Dynamics (CFD) problems. By leveraging SuperNODE's unprecedented capability to support up to 32 GPUs and directly utilize terabytes of GPU memory, engineers and researchers can solve and visualize massive CFD simulations efficiently and at unprecedented speeds. This innovation not only revolutionizes the field of CFD by drastically reducing the time and resources required for complex simulations, but also demonstrates a remarkable synergy of high-performance computing and sophisticated software development, paving the way for new horizons in various industries that rely on fluid dynamics analysis.

"FluidX3D and GigaIO SuperNODE are a perfect match for solving extreme-scale CFD problems. This open, flexible solution allows me to use terabytes of VRAM directly on the maximum number of GPUs. It significantly accelerates computing and enables me to solve and render some of the largest and most challenging CFD problems, intractable on alternate hardware platforms."

– Dr. Moritz Lehmann, Sole Developer of FluidX3D

