

Transforming High-Performance Computing with VAST Data



Table of Contents

Chapter 1: Introduction to High-Performance Computing	5
What is HPC?	5
Why Use HPC?	5
Understanding Supercomputers and HPC Clusters	6
The Role of GPU-Based HPC	6
HPC Terminology and Lingo	6
How VAST Accelerates HPC Workflows	8
Chapter 2: Challenges in HPC Today	9
Addressing HPC Growth Needs	9
Performance Bottlenecks	9
Data Movement and Data Gravity Problem	10
Scalability Challenges	10
Complexity in Data Management and Administration	11
The Need for Real-Time Data Processing in HPC and AI	12
Summary of Key HPC Challenges	12
Chapter 3: HPC with VAST Data: Key Differentiators	13
Transforming HPC Infrastructure	13
Disaggregated Shared-Everything Architecture	13
All-Flash Storage Optimized for High Throughput	14
Integration with NVIDIA BlueField DPUs and Enhanced AI Capabilities	14
Multiprotocol Support with No Data Translation	15
Enhanced Data Management and Simplified Administration	16
Summary of VAST Data's Key Differentiators for HPC	17

Chapter 4: HPC with VAST Data Application Examples by Market	18
VAST Data's Transformative Impact Across Industries	18
Healthcare and Life Sciences	18
Medical Imaging and Diagnostics	19
Higher Education and Research	19
Financial Services	20
Media and Entertainment	20
Automotive	21
Cloud Service Providers	21
Energy and Utilities	21
Government	22
Telecommunications	22
Summary of Key HPC with VAST Data Application Examples by Market	23
Chapter 5: Summary and Key Takeaways	24
Evolving HPC with VAST Data	24
The Role of HPC and VAST Data's Differentiators	24
VAST Data's Impact Across Key Industries	25
Simplified Administration and High Availability	26
Future-Ready Storage for AI-Driven HPC	26
The Power of Unified HPC and AI Workloads	27
Key Takeaways for VAST Data's Value in HPC	27
Conclusion	28

Chapter 1: Introduction to High-Performance Computing

What is HPC?

High-Performance Computing (HPC) refers to the use of powerful computational resources and complex architectures to solve demanding scientific, engineering, and business problems at unprecedented speeds. Unlike conventional computing, HPC systems can perform trillions of calculations per second, enabling researchers and organizations to tackle large-scale simulations, complex data analyses, and artificial intelligence tasks that were once beyond reach.

HPC systems typically consist of clusters of high-speed processors and specialized interconnects, which allow for parallel processing. In this method, multiple processors simultaneously work on different parts of a task. This parallelism is at the heart of HPC, enabling high-throughput, low-latency operations across large datasets or computational models.

Refer to the “[Modernizing HPC Infrastructure and Administration with VAST Data](#)” white paper for additional insights into HPC architectures.

Why Use HPC?

HPC is essential across industries where the speed and scale of computation can lead to breakthrough discoveries, faster insights, and significant efficiency gains. Some critical applications include:

- **Scientific Research:** Fields such as genomics, climate modeling, and quantum mechanics rely on HPC to perform simulations and analyze vast amounts of data.
- **Healthcare:** HPC aids in drug discovery, genomics research, and precision medicine by analyzing large-scale datasets and modeling complex biological processes.
- **Financial Services:** Risk modeling, fraud detection, and high-frequency trading benefit from HPC’s ability to rapidly handle complex calculations.
- **Engineering and Manufacturing:** Companies use HPC for simulations in areas such as aerodynamics, structural analysis, and materials science.

By reducing computational times and handling high-resolution data, HPC systems enable industries to achieve outcomes that would be impossible with conventional computing resources.

Understanding Supercomputers and HPC Clusters

Supercomputers are highly specialized machines engineered for large-scale computations at unmatched speeds. These systems comprise thousands—or even millions—of interconnected nodes (or processors) that work in parallel to perform calculations. Their performance is measured in FLOPS (Floating Point Operations Per Second), with the world's most advanced systems achieving exaflop performance—capable of executing a quintillion calculations per second. Supercomputers are often custom-built to address highly specific computational challenges like climate modeling, quantum physics, or genomic research.

In contrast, HPC clusters are collections of standardized computing servers or nodes connected via high-speed networks to function as a unified system. While they aim to enable large-scale computations, HPC clusters differ in their design and approach from supercomputers. They are typically built using commercially available, off-the-shelf hardware, making them a more flexible, cost-effective, and scalable solution for a wide range of scientific, industrial, and academic applications.

Modern HPC clusters, as detailed in the [“Modernizing HPC Infrastructure and Administration with VAST Data”](#) whitepaper, enable incremental deployment of computing power, combining cost-effectiveness with scalability for workloads of varying complexities.

The Role of GPU-Based HPC

In recent years, Graphics Processing Units (GPUs) have transformed HPC by enabling even faster parallel processing capabilities than conventional CPU-based systems. Initially designed for rendering graphics, GPUs excel at handling complex mathematical calculations across numerous threads, making them ideal for HPC workloads.

GPU-based HPC systems are integral to fields like Artificial Intelligence (AI) and deep learning, where massive amounts of data must be processed in parallel. These systems allow researchers to train large machine learning models and run advanced simulations with speed and efficiency. As a result, GPU-based HPC has accelerated the pace of discovery and innovation in various fields, enabling breakthroughs previously thought to be impossible.

HPC Terminology and Lingo

To help navigate the specialized language of HPC, here are some common terms and their definitions:

- **Node:** A single computer within an HPC cluster. Nodes are typically classified as compute nodes (responsible for processing tasks), head nodes (serving as the primary interface for users to submit jobs, access data, and manage workflows), or management nodes (overseeing cluster operations, including job scheduling, resource allocation, monitoring, and system provisioning).

- **Cluster:** A collection of interconnected nodes functioning as a single computational resource.
- **Parallel Processing:** A computing technique where multiple processors simultaneously handle different parts of a task, essential for HPC performance.
- **Job Scheduling:** The process of allocating computing resources to different tasks or “jobs” in an HPC environment. Popular job schedulers include Slurm, Portable Batch System (PBS), and Load Sharing Facility (LSF).
- **FLOPS:** Floating Point Operations Per Second, a measure of computational power. HPC systems are typically measured in teraflops, petaflops, or exaflops.
- **Data Gravity:** The concept that large datasets “pull” applications and services toward them due to the cost and complexity of moving data, a significant factor in HPC and data storage.
- **Network Bandwidth:** This is a network’s capacity to transmit data. In HPC, high-bandwidth networks are necessary to ensure fast data transfer between nodes.
- **Latency:** The delay in data transfer within an HPC system. Low latency is critical for high-speed communication between processors in HPC.
- **NVMe-oF (NVMe over Fabrics):** A protocol enabling high-speed access to remote NVMe storage over network fabrics like Ethernet (RoCE) and InfiniBand, minimizing latency for demanding workloads.
- **RDMA (Remote Direct Memory Access):** A technology that allows direct memory-to-memory data transfer between systems, reducing latency and offloading the CPU for better performance in HPC.
- **RoCE (RDMA over Converged Ethernet):** A protocol that delivers RDMA functionality over Ethernet networks, improving data transfer efficiency in HPC clusters.

Understanding these terms is essential for anyone engaging with HPC, as they frequently appear in discussions about system performance, architecture, and functionality.

How VAST Accelerates HPC Workflows

Imagine a research team analyzing genomic data to develop new treatments. Every day, they process petabytes of data from thousands of genetic sequences. Speed is critical—they need results fast to move the project forward. This is where the VAST Data Platform, with NVMe-oF (NVMe over Fabrics), RDMA (Remote Direct Memory Access), and RoCE (RDMA over Converged Ethernet), becomes the backbone of their operation.

When the scientists submit a job, the compute nodes spring into action, contacting the VAST Data storage system. Instead of waiting for the slow read times typical of conventional storage, NVMe-oF enables these nodes to access massive data volumes at speeds comparable to having the data locally stored on the nodes themselves.

Meanwhile, RDMA allows the compute nodes to pull data directly from the storage, bypassing CPU bottlenecks. This direct memory access ensures the data flows smoothly without interruption, significantly reducing processing delays. It's like having a dedicated express lane for data transferring between storage and compute—no stopping, no slowing down. This scenario occurs over RoCE, which carries RDMA over the Ethernet network.

With the VAST Data Platform running NVMe-oF, RDMA, and RoCE, researchers experience ultra-fast data processing, completing tasks that would typically take days in just hours. As data flows effortlessly from storage to compute, insights come to life faster, accelerating breakthroughs and keeping the research moving forward at lightning speed.

Chapter 2: Challenges in HPC Today

Addressing HPC Growth Needs

HPC systems are essential in addressing complex problems across industries, including genomics, hedge funds, climate science, engineering simulations, and AI. However, as data volumes grow and computation becomes more intensive, HPC environments face several persistent challenges. These include performance bottlenecks, data movement inefficiencies, scalability limitations, and the need for high-availability architectures to support enterprise-level operations. Here, we'll explore these challenges and introduce how modern solutions, like the VAST Data Platform, can help address these pain points.

Performance Bottlenecks

HPC environments are built to handle massive datasets and execute high-throughput computations. Yet, they often face performance bottlenecks due to the limitations of conventional storage and data access methods.

Challenges in Data Access and Throughput

Conventional parallel file systems like Lustre, and GPFS rely on outdated architectures initially designed for the hard drive era. These systems often struggle with the high throughput demands of modern workloads, particularly those involving AI and machine learning, where high-speed data access is critical. Data bottlenecks can occur when multiple compute nodes attempt to access the same storage resources simultaneously, causing latency and reducing overall system performance.

Storage I/O Limitations

In HPC environments, storage I/O (input/output) speeds play a critical role in achieving desired computational speeds. Many conventional HPC storage solutions rely on spinning disks or hybrid architectures that cannot fully leverage the speed advantages of solid-state drives (SSDs). This can result in slower data retrieval and storage access times, impeding computational workflows and forcing organizations to look for alternative, more efficient data architectures.

How VAST Data Solves These Bottlenecks

VAST Data's all-flash architecture, combined with its Disaggregated Shared-Everything (DASE) model, maximizes data throughput and minimizes latency, enabling HPC environments to achieve the necessary data speeds. By allowing data to be accessed across the Global Namespace without the need for data movement, VAST Data effectively eliminates conventional I/O bottlenecks, ensuring compute nodes can access data with minimal delay.

Data Movement and Data Gravity Problem

The concept of “data gravity” describes the phenomenon where large datasets naturally pull applications and services closer to them, as moving data can be costly and complex. In HPC environments, data gravity creates significant challenges, especially when datasets need to be processed in multiple locations or shared across different systems.

Conventional Data Movement Challenges

Conventional HPC systems often require data to be moved between storage locations, compute nodes, or even physical sites to perform specific tasks. This constant movement not only consumes time and bandwidth, but also introduces risks related to data integrity, security, and latency. These issues become more pronounced as data volumes grow, creating inefficiencies that slow down workflows and increase costs.

Impact on Multi-Location and Cloud-Connected HPC Environments

Data movement between locations can create operational challenges in environments spanning multiple sites, such as collaborative research institutions or geographically distributed HPC clusters. For example, transferring large datasets to a cloud provider or across regions may require high-bandwidth connections and specialized data transfer solutions, adding complexity and cost.

How VAST Data Addresses the Data Gravity Problem

VAST Data’s DASE architecture allows data to remain in place while being accessible across different locations. This approach eliminates the need for pre-movement, reducing complexity, costs, and latency. With a global namespace, data remains available across edge, core, and cloud environments, simplifying the collaboration and processing of large datasets without extensive data movement.

Scalability Challenges

HPC environments must scale efficiently to accommodate growing workloads and evolving computational needs. Conventional storage systems, however, often struggle to scale without introducing complexity, high costs, and management overhead.

Hardware and Software Limitations in Scaling

Conventional parallel file systems like Lustre and GPFS often require extensive tuning and complex configurations to scale effectively. Adding new nodes or expanding storage in these environments can require reconfiguring file striping, adjusting metadata servers, and carefully managing storage pools, all of which add administrative complexity and increase the risk of downtime.

Impact of Increasing Data Volume and Velocity

With the rise of AI, big data, and machine learning, the volume and velocity of data are increasing rapidly. Many conventional HPC systems were not designed to handle the exponential data growth, resulting in performance degradation as more data is added to the system. This is particularly problematic in research and enterprise HPC environments, where data access speeds must keep up with large-scale computational workloads.

How VAST Data Eases Scaling Challenges

VAST Data provides a highly scalable architecture that allows organizations to add storage capacity without reconfiguring or tuning. The VAST Data Platform is designed for high-throughput, low-latency operations that maintain performance regardless of data volume, making it ideal for environments where scalability and simplicity are essential. VAST Data enables seamless scaling across edge, core, and cloud environments by eliminating data islands and offering a unified data architecture.

Complexity in Data Management and Administration

Managing HPC systems can be challenging, as these environments require specialized knowledge and resources to handle configurations, data transfers, security, and system tuning.

Administrative Complexity in Conventional HPC Systems

Parallel file systems such as Lustre, and GPFS require administrators to configure various settings, including data striping, replication, and caching policies. Each of these configurations impacts performance, stability, and data redundancy. As HPC environments grow, managing these settings becomes more complex, requiring significant expertise and manual intervention to optimize systems.

Challenges in Maintaining Data Consistency and Reliability

In addition to complexity, conventional HPC systems often struggle to maintain data consistency across nodes, particularly in multi-site environments. These systems can be prone to data corruption, inconsistent snapshots, and other reliability issues that compromise data integrity, necessitating additional tools and resources.

How VAST Data Simplifies Data Management

VAST Data's architecture consolidates storage, management, and computing, reducing administrative complexity and simplifying data management. With features like automated data placement, snapshots, and deduplication, the VAST Data Platform minimizes the need for manual configurations and third-party tools, making it easier for administrators to maintain data integrity. VAST Data's seamless failover capabilities also ensure high availability, reducing the risk of downtime and data loss.

Further operational insights are detailed in the "[Modernizing HPC Infrastructure and Administration with VAST Data](#)" whitepaper.

The Need for Real-Time Data Processing in HPC and AI

As HPC and AI applications continue to converge, there is an increasing demand for real-time data processing.

Conventional HPC systems are often optimized for batch processing, which can introduce delays in data availability and limit the effectiveness of AI and machine learning workloads that depend on immediate insights.

Batch Processing Limitations in Conventional HPC Systems

Most conventional HPC environments are designed for batch processing, where tasks are queued and processed sequentially. This approach can lead to delays, as new data must wait in line to be processed, reducing system responsiveness. For AI applications, where real-time analysis and model training are critical, batch processing is often insufficient.

Challenges in Managing Data for AI Workloads

AI workloads require continuous data feeds and low-latency data access. However, integrating AI with conventional HPC systems often requires significant configuration and management, as these systems were not designed to support the constant data retrieval and processing that AI demands.

How VAST Data Enables Real-Time Data Processing

VAST Data's InsightEngine, developed in collaboration with NVIDIA, transforms data processing by enabling real-time ingestion, vector embeddings, and semantic data processing. Unlike batch-based systems, InsightEngine processes incoming data immediately, supporting Retrieval-Augmented Generation (RAG) and AI-driven applications that benefit from real-time insights. VAST Data's scalable architecture also allows for trillions of vector embeddings, enabling extensive real-time similarity searches and accelerated data analysis.

Summary of Key HPC Challenges

Today's HPC environments face a set of persistent challenges, including:

1. **Performance Bottlenecks:** Conventional storage and I/O limitations impede computational speeds.
2. **Data Movement and Data Gravity:** The high cost and complexity of moving large datasets can hinder workflows.
3. **Scalability Challenges:** Conventional systems struggle to scale seamlessly, often requiring complex reconfigurations.
4. **Administrative Complexity:** Managing data across nodes and sites demands specialized expertise and resources.
5. **The Need for Real-Time Processing:** Batch processing limitations slow down AI and machine learning applications that rely on immediate data insights.

These challenges underscore the need for a modern HPC platform capable of delivering real-time data processing, seamless scalability, and high performance. In the next chapter, we'll explore the VAST Data Platform's key differentiators. See the "[Modernizing HPC Infrastructure and Administration with VAST Data](#)" whitepaper for additional information.

Chapter 3: HPC with VAST Data: Key Differentiators

Transforming HPC Infrastructure

The VAST Data Platform introduces a transformative approach to handling the scale, complexity, and demands of modern HPC environments. Unlike conventional storage solutions designed for the hard-drive era, VAST Data's all-flash, DASE architecture delivers high-throughput, low-latency access with unparalleled simplicity and efficiency.

This chapter will explore the differentiators that set VAST Data apart from conventional HPC storage solutions, focusing on features that reduce data movement, optimize performance, and simplify administration.

Disaggregated Shared-Everything Architecture

The VAST Disaggregated Shared-Everything (DASE) architecture represents a fundamental shift from conventional HPC storage systems, which often rely on a shared-nothing model. In a shared-nothing architecture, each storage node operates independently, limiting scalability and creating bottlenecks. By contrast, VAST's DASE model decouples compute from storage, allowing storage resources to be pooled and accessed by any compute resource across the cluster.

Benefits of Disaggregated Shared-Everything

- **Unlimited Scalability:** This architecture enables organizations to add storage or compute resources independently, providing true scalability without the limitations of conventional storage clusters.
- **Unified Data Access:** Data is accessible across edge, core, and cloud environments, allowing for seamless data processing and eliminating the need to move data between locations for different processing phases.
- **Improved Resilience and Performance:** By removing the dependence on shared-nothing components, VAST Data's architecture ensures high availability and minimizes single points of failure, reducing downtime and increasing reliability.

Addressing the Data Gravity Problem

VAST Data's architecture solves the data gravity problem by providing a global namespace that spans locations, allowing data access without requiring physical movement. This flexibility allows organizations to consolidate data without creating isolated silos, making it easier to collaborate on large datasets and reducing data transfer costs. With the VAST Data Platform, data remains where it's needed, making it accessible across geographically distributed teams and resources.

All-Flash Storage Optimized for High Throughput

Conventional HPC storage solutions often rely on a mix of hard drives and SSDs, limiting their ability to fully leverage flash performance. In contrast, VAST Data's all-flash infrastructure delivers the high-speed data access modern workloads require.

High-Performance NVMe-Based Storage

The VAST Data Platform is built on NVMe-based storage, enabling ultra-low latency and high throughput for both small and large files. This architecture is ideal for HPC applications that demand rapid data access, such as AI training, genomics, hedge funds, and climate modeling. By leveraging the latest in NVMe technology, VAST Data eliminates the I/O bottlenecks that conventional systems experience, ensuring data access speeds are no longer a limiting factor in computation.

Consistent Performance Across Workloads

VAST Data's all-flash architecture guarantees consistent, high-performance data access across a variety of workloads, whether they involve large sequential reads or numerous small, random I/O requests. This makes the VAST Data Platform an excellent choice for mixed-use environments where multiple applications with varying data demands run simultaneously, ensuring every task gets the data it needs without delays.

Real-World Impact on HPC Workflows

For HPC environments, VAST Data's high-throughput architecture results in faster processing times, reduced wait times for data access, and more efficient resource utilization. This high performance accelerates computational tasks and supports the increasingly real-time needs of AI and machine learning applications.

Integration with NVIDIA BlueField DPUs and Enhanced AI Capabilities

As HPC converges with AI, the demand for low-latency, real-time data access continues to grow. The VAST Data Platform is optimized for NVIDIA BlueField DPUs (Data Processing Units), bringing the compute and storage closer together to enable efficient HPC and AI workloads.

Offloading Compute with NVIDIA BlueField DPUs

By integrating with NVIDIA's DPUs, VAST Data offloads storage-related tasks from the CPU to the DPU. This integration moves the storage controller functions directly onto the GPU server, freeing up CPU cycles for more intensive computing tasks. Unlike conventional storage solutions, where storage operations can slow down compute performance, VAST's DPU integration enables efficient processing without compromising computational resources.

Supporting AI Workloads and Real-Time Data Processing

VAST Data's collaboration with NVIDIA enhances its platform's support for AI-driven HPC workloads. Real-time data processing, powered by the VAST InsightEngine with NVIDIA, allows AI applications to access data instantaneously, transforming incoming data into vector embeddings and graph relationships as soon as it is ingested. This real-time processing supports RAG and complex AI models that require immediate data access, making the VAST Data Platform ideal for modern HPC/AI integrations.

Multiprotocol Support with No Data Translation

In HPC environments, data access requirements vary significantly across workflows, making multiprotocol support essential for efficient storage use. While many vendors claim to offer multiprotocol support by enabling access to parallel file system data through NFS or SMB, these solutions typically rely on gateways that translate between protocols. This approach has significant drawbacks, including reduced speed, limited scalability, and potential data inconsistency.

VAST Data's Multiprotocol Advantage: Direct Access without Gateways or Data Translation

VAST Data's DASE architecture eliminates the need for gateways by providing true multiprotocol support. This allows NFS, SMB, S3, and Block clients to access the same data directly without requiring translation or creating additional data copies. Here's why this approach is superior:

- 1. No Data Translation or Copies:** Competitor solutions often require translation processes or multiple data copies to enable cross-protocol access. This can be inefficient, leading to delays and increased complexity. VAST Data's direct access capability avoids this by allowing all protocols to read and write directly to a single, unified data layer. This eliminates the overhead of translation and reduces the risk of data inconsistencies.
- 2. Eliminating Gateways for Faster Performance:** Many vendors use gateways to enable NFS or SMB access to parallel file systems. These gateways create bottlenecks that can reduce performance by up to 90% compared to direct access. With VAST Data, there is no gateway layer, which means NFS, SMB, S3, and Block clients can interact directly with the data at native NVMe speeds, providing up to 10 times faster performance.
- 3. Scalability without Bottlenecks:** Gateway architectures are typically not designed to scale at the same level as native access methods, often limiting scalability to a fraction of what direct access can achieve. VAST Data's multiprotocol support scales seamlessly without gateways, enabling HPC environments to handle high-throughput workloads and large user bases without compromising on performance or scalability.

- 4. Simplified Data Management:** With VAST Data's direct multiprotocol support, all clients access the same data layer regardless of protocol. This unified access eliminates the need to maintain separate data copies for different protocols, reducing storage overhead and simplifying data management. For HPC environments with mixed workloads, this means that AI, analytics, and simulation tasks can all access data seamlessly, simplifying workflows and improving productivity.

VAST Data Advantage: True Multiprotocol Support without Compromise

VAST Data's ability to provide NFS, SMB, S3, and Block access directly to the same dataset, without translation or gateways, sets it apart from competitors. This approach not only maximizes performance and scalability but also reduces operational complexity. By eliminating data translation and providing unified access across protocols, the VAST Data Platform delivers a powerful solution that enhances flexibility, simplifies management, and enables high-speed data access across the HPC environment.

Enhanced Data Management and Simplified Administration

One of the primary challenges in conventional HPC environments is the complexity of configuring and managing parallel file systems. VAST's approach significantly reduces administrative overhead, enabling efficient data management and access without the need for extensive tuning.

Automated Data Placement and Management

VAST automates many aspects of data management, including data placement, caching, and load balancing. This automation allows administrators to focus on high-value tasks rather than manual configurations, reducing the risk of human error. Overall, it ensures optimized performance without the need for constant adjustments.

Integrated Snapshots and Deduplication

The VAST Data Platform includes snapshots and deduplication capabilities directly in its architecture. Conventional HPC storage systems often lack coherent snapshot functionality or require third-party solutions for deduplication, increasing costs and complexity. VAST Data's integrated approach allows administrators to create space-efficient snapshots and seamlessly eliminate redundant data, simplifying data management while reducing storage requirements.

High Availability and Reliability

VAST Data provides built-in High Availability (HA) and failover capabilities, ensuring data accessibility even in the event of hardware failure. This is particularly important in HPC environments where downtime can disrupt critical workflows. Unlike Lustre or GPFS, which require third-party tools for HA, the VAST Data Platform's architecture is designed to deliver continuous access, minimizing potential disruptions and ensuring operational reliability.

Summary of VAST Data's Key Differentiators for HPC

VAST Data brings a new level of flexibility, performance, and simplicity to HPC environments through its DASE architecture and all-flash infrastructure. Key differentiators include:

1. **DASE Architecture:** Enables scalability and reduces data movement by providing a global namespace that spans multiple locations.
2. **High-Performance All-Flash Storage:** NVMe-based architecture ensures low-latency, high-throughput data access, suitable for intensive HPC and AI workloads.
3. **NVIDIA BlueField DPUs Integration:** Offloads compute tasks to DPUs to enhance computational efficiency without slowing down performance.
4. **Multiprotocol Support with No Data Translation:** Allows seamless data access across NFS, SMB, S3, and Block without data duplication or conversion.
5. **Simplified Data Management:** Automated data placement, snapshots, deduplication, and high availability reduce the need for manual configurations and administrative effort.
6. **InsightEngine with NVIDIA for Real-Time AI and Data Processing:** Provides real-time data ingestion, vector embeddings, and scalability for AI applications, consolidating storage and processing in a single system.

These capabilities make the VAST Data Platform an ideal choice for modern HPC environments, enabling organizations to overcome the limitations of conventional storage systems and meet the demands of a data-intensive future. For comparisons with Lustre and GPFS, refer to the [“Modernizing HPC Infrastructure and Administration with VAST Data”](#) whitepaper.

Chapter 4: HPC with VAST Data

Application Examples by Market

VAST Data's Transformative Impact Across Industries

The VAST Data Platform's revolutionary Disaggregated Shared-Everything (DASE) architecture delivers unmatched scalability, high-throughput performance, and simplified data access, empowering diverse industries to tackle the complexities of HPC and data-intensive workflows. By eliminating bottlenecks and enabling real-time access to critical data, VAST helps organizations unlock new levels of productivity and accelerate innovation across a wide range of applications.

From genomics research and AI-driven diagnostics in Healthcare and Life Sciences to climate modeling and cross-institutional collaboration in Higher Education and Research, VAST empowers researchers and practitioners to process vast datasets and gain actionable insights faster than ever. In the Financial Services and Telecommunications industries, VAST's low-latency storage architecture supports split-second decision-making and fraud detection, while Energy and Utility companies leverage VAST for seismic data analysis and smart grid optimization. Meanwhile, the Media & Entertainment and Automotive sectors benefit from VAST's high-performance data access, enabling seamless video editing, visual effects rendering, and autonomous vehicle development.

This chapter provides an in-depth look at key use cases across these sectors, demonstrating how VAST Data addresses the unique challenges of each industry. By supporting diverse, high-demand applications with a single, unified platform, VAST helps organizations not only manage current data needs but also scale for the future, fueling continuous innovation and industry leadership.

Healthcare and Life Sciences

The Healthcare and Life Sciences sector depends on extensive data-driven applications in genomics, drug discovery, and precision medicine. VAST Data's architecture, designed for high throughput and scalability, meets the sector's data demands, driving innovation in patient care and research.

- **Precision Medicine and Genomics Research:** Genomics workflows require rapid, high-throughput data access for analyzing genetic sequences. VAST Data's all-flash storage enables real-time access without data pre-movement, accelerating primary, secondary, and tertiary analyses in genomics. VAST's InsightEngine with NVIDIA further supports complex data processing by enabling RAG, allowing researchers to identify patterns across multiple diseases.
 - **Customer Perspective:** In pediatric disease research, VAST Data facilitates efficient genome processing, reducing time-to-insight, which is crucial for developing targeted treatments.

- **Drug Discovery and Translational Research:** Drug discovery requires extensive analysis of molecular structures and patient data, and VAST's high-performance data access shortens this process. VAST's DASE architecture allows real-time access to drug datasets and reduces the need for data duplication, accelerating AI-based discovery and model training.
 - **Customer Perspective:** Research teams using VAST report faster candidate identification and a streamlined discovery pipeline, allowing them to bring potential treatments to trial more quickly.

Medical Imaging and Diagnostics

Medical imaging in healthcare produces massive datasets that require efficient storage and retrieval. VAST's architecture ensures seamless data access, improving diagnostic efficiency and enabling advanced AI-driven imaging solutions.

- **AI-Enhanced Imaging:** VAST's real-time access supports machine learning in imaging, allowing algorithms to process and analyze medical images quickly for anomaly detection. By integrating InsightEngine with NVIDIA, radiology teams can use AI to interpret data more efficiently.
 - **Customer Perspective:** Hospitals using VAST's solutions report quicker image processing and retrieval, enhancing patient care with faster diagnostic workflows and improved radiologist productivity.

Higher Education and Research

Higher education and research institutions rely on VAST's scalable, high-performance platform to manage large datasets, facilitating collaborative and interdisciplinary research.

- **Multidisciplinary Research:** VAST's global namespace eliminates data silos, allowing research departments to collaborate more effectively. Our multiprotocol support ensures that researchers from different fields can access shared data seamlessly, fostering innovation across departments like climate science, physics, and engineering.
 - **Customer Perspective:** Universities using VAST report accelerated data sharing between departments, improving collaborative research and increasing efficiency in cross-disciplinary projects.
- **Cross-Institutional Data Sharing:** With VAST's DASE architecture, researchers can access data across multiple locations in real-time, supporting large-scale, multi-institutional studies.
 - **Customer Perspective:** The VAST Data Platform has enabled universities to participate in global research projects with fewer delays, expanding collaboration capabilities and improving research outcomes.

Financial Services

In financial services, applications such as trading, fraud detection, and risk analysis depend on real-time data processing. VAST's architecture supports these high-demand tasks by reducing latency and providing reliable data access.

- **Algorithmic Trading and Risk Analysis:** Low-latency access is essential in trading environments. The VAST Data Platform provides high-throughput, consistent data access, enabling financial institutions to execute trades and analyze risks in real-time. Our DASE architecture reduces data retrieval times, improving trading models' response capabilities.
 - **Customer Perspective:** Financial firms using VAST report enhanced market data access, enabling rapid risk analysis and trade execution, essential in fast-paced financial markets.
- **Fraud Detection and Compliance:** VAST's real-time processing and global namespace allow fraud detection algorithms to analyze transaction patterns at scale, accelerating detection and reducing potential losses.
 - **Customer Perspective:** Financial firms report faster fraud pattern recognition, enhancing overall security and compliance.

Media and Entertainment

Media production requires high-throughput storage for rapid access to large video files and visual effects assets. VAST's architecture provides the speed and reliability required for time-sensitive production schedules.

- **Video Editing and Production:** VAST's high-throughput, all-flash storage ensures smooth access to video files for real-time editing. This access allows production teams to work seamlessly, reducing bottlenecks during post-production.
 - **Customer Perspective:** Media companies using VAST report improved production timelines, with editors experiencing minimal wait times when loading large files.
- **Animation and VFX:** VAST's direct data access enables quick retrieval of high-resolution assets, supporting animation and VFX teams in real-time rendering and accelerating the iterative creative process.
 - **Customer Perspective:** Animation studios leveraging VAST have reduced delays in asset access, enabling higher productivity and faster time-to-completion.

Automotive

In the automotive industry, VAST supports data-intensive applications, such as AI model training for autonomous vehicles and large-scale simulations, both of which are crucial to vehicle development.

- **Autonomous Driving Development:** Autonomous vehicles rely on large datasets for training AI models. The VAST Data Platform accelerates this process with real-time data access, allowing automotive engineers to develop, test, and iterate models faster.
 - **Customer Perspective:** Automotive R&D teams report shorter AI model training times, driving faster advancements in autonomous technology.
- **Simulation and Testing:** VAST's DASE architecture provides consistent access to simulation data, enabling automotive engineers to run tests on vehicle performance, safety, and aerodynamics with minimal delays.
 - **Customer Perspective:** Automotive engineers using VAST experience faster access to simulation data, enabling quicker testing and design improvements, accelerating time-to-market for innovations.

Cloud Service Providers

Cloud Service Providers (CSPs) require scalable, high-performance storage that can support diverse workloads in multi-tenant environments. VAST's architecture delivers reliable data access and efficient scaling, meeting the high-performance requirements of cloud providers.

- **Multi-Tenant Workload Support:** The DASE architecture provides scalable, secure storage that meets the demands of multi-tenant environments, enabling efficient workload management and consistent performance across customer applications.
 - **Customer Perspective:** Cloud Service Providers using VAST report higher reliability in workload distribution, ensuring performance stability across diverse applications.

Energy and Utilities

Energy and utility companies rely on high-throughput storage to analyze the vast datasets, such as seismic or grid data, needed for resource management and infrastructure monitoring.

- **Seismic Data Analysis for Oil and Gas:** VAST's platform supports real-time analysis of seismic data, helping companies explore and model geological structures faster.
 - **Customer Perspective:** Oil and gas companies using VAST report faster resource discovery and more efficient geological analysis, reducing project times and costs.

- **Smart Grid Monitoring and Data Analytics:** For utilities, VAST's DASE architecture supports real-time grid data analytics, enabling efficient resource management and proactive infrastructure maintenance.
 - **Customer Perspective:** Utility providers using VAST report faster grid analysis, allowing proactive adjustments to ensure stability and reduce costs.

Government

Government agencies rely on HPC to support critical applications such as disaster response simulations and defense research, where timely data access and scalability are paramount.

- **Disaster Response Modeling:** The VAST Data Platform enables real-time analysis for disaster simulations, helping agencies respond more effectively to crisis scenarios. With VAST's multiprotocol support, agencies can collaborate seamlessly across departments.
 - **Customer Perspective:** Agencies leveraging VAST report improved coordination and faster data analysis during disaster simulations, enhancing their ability to prepare for and manage emergencies.
- **Defense Research:** Real-time data access is crucial for secure, efficient analysis in defense applications. VAST's architecture provides the reliability and performance necessary for high-stakes government operations.
 - **Customer Perspective:** Defense organizations using VAST experience improved data processing speeds and enhanced system reliability when supporting mission-critical applications.

Telecommunications

Telecommunications companies require fast data access for network analysis, customer experience optimization, and AI-driven applications. VAST's platform provides the speed and scalability necessary for these high-demand applications.

- **Network Optimization:** VAST's DASE architecture allows telecom providers to analyze large datasets in real-time, enabling quicker adjustments to optimize network performance and improve service quality.
 - **Customer Perspective:** Telecom companies report improved response times to network changes, resulting in enhanced service delivery.

Summary of Key HPC with VAST Data Application Examples by Market

This chapter showcased how VAST Data's unique architecture supports diverse HPC and data-intensive applications across industries, ranging from Healthcare and Life Sciences to Telecommunications and Automotive. By eliminating conventional bottlenecks in data movement, storage access, and scalability, VAST enables faster data processing, real-time insights, and streamlined collaboration, transforming how organizations approach complex data challenges.

Through case studies and examples, we explored how VAST's DASE architecture delivers critical benefits across sectors: empowering researchers in genomics, accelerating algorithmic trading and risk analysis in financial services, enhancing simulation capabilities in higher education, and optimizing network performance in telecommunications. Each example demonstrates VAST's ability to meet industry-specific needs while providing a flexible, scalable platform that evolves alongside the organization.

As industries continue to push the boundaries of HPC, the VAST Data Platform emerges as a future-proof solution, empowering organizations to drive innovation, boost productivity, and achieve results faster. Refer to the "[Modernizing HPC Infrastructure and Administration with VAST Data](#)" whitepaper for additional discussion on how the VAST Data Platform is transforming HPC environments worldwide.

Chapter 5: Summary and Key Takeaways

Evolving HPC with VAST Data

The landscape of HPC is transforming, driven by growing demands for fast, reliable, and scalable data access across research and industry applications. The VAST Data Platform meets these challenges head-on, providing solutions that overcome the limitations of conventional HPC storage systems. By introducing a DASE architecture and leveraging NVMe-oF and RDMA, VAST Data empowers organizations to elevate their computational capabilities, reduce complexity, and enhance scalability.

The Role of HPC and VAST Data's Differentiators

HPC is fundamental for tackling complex problems in scientific research, healthcare, engineering, and AI. However, conventional HPC storage solutions, such as Lustre, and GPFS, face significant constraints in scalability, performance, and ease of management. VAST Data eliminates these issues by employing a DASE architecture, enabling the real-time data processing that is critical for AI-integrated HPC workloads. Key differentiators include:

- **The DASE Architecture:** Decouples computing from storage, providing scalable, flexible data access without the constraints of conventional shared-nothing systems.
- **All-Flash NVMe-Based Storage:** Delivers high throughput and low latency, ideal for both high-sequential and high-random I/O workloads.
- **Integration with NVIDIA BlueField DPUs:** Optimizes compute resources by offloading storage tasks to DPUs, enabling real-time data processing and AI workflows.
- **Multiprotocol Support:** Enables seamless access across NFS, SMB, S3, and Block without data duplication, simplifying workflows and improving application interoperability.
- **InsightEngine with NVIDIA for Real-Time AI:** Facilitates real-time data ingestion and processing, supporting AI-driven HPC tasks and RAG.

These features collectively position VAST Data as an ideal solution for HPC environments requiring high performance, scalability, and ease of management.

VAST Data's Impact Across Key Industries

The VAST Data Platform supports a wide array of industries, helping organizations manage complex HPC and AI-driven workflows through high-speed, low-latency data access. This transformative impact spans critical sectors, enabling innovation, efficiency, and real-time insights across diverse applications.

- **Life Sciences and Healthcare:** In genomics research and precision medicine, VAST accelerates workflows by providing rapid, high-throughput data access without the need for data movement between storage tiers. Researchers benefit from real-time processing and analysis, crucial for fields like pediatric cancer research, where swift data processing is essential for treatment development. VAST's DASE architecture enables smooth data retrieval and enhanced research capabilities, directly advancing medical research and patient care.
- **Medical Imaging and Diagnostics:** VAST's architecture optimizes medical imaging processes, supporting AI-driven diagnostics by providing immediate data access to large imaging datasets. By using VAST InsightEngine with NVIDIA, medical facilities can leverage machine learning and analyze images faster, enhancing diagnostic workflows and improving patient outcomes.
- **Higher Education and Research:** VAST supports large-scale, interdisciplinary research by simplifying data sharing and enabling seamless collaboration across departments and institutions. Its global namespace eliminates data silos, allowing researchers in fields like climate science and engineering to easily access shared data, fostering cross-departmental research, and supporting multi-institutional projects.
- **Financial Services:** VAST enables real-time data access for algorithmic trading, fraud detection, and compliance. Its low-latency storage architecture enhances transaction speed and risk analysis by providing consistent, high-throughput access to market data, helping institutions make faster, more accurate data-driven decisions.
- **Media and Entertainment:** VAST supports high-throughput data access for video production, visual effects, and animation. With real-time data access, media professionals can edit, render, and process video files efficiently, reducing bottlenecks and enhancing productivity in post-production workflows.
- **Automotive:** In the automotive industry, VAST enables faster simulation and AI model training for autonomous driving development. By ensuring low-latency access to simulation and sensor data, VAST supports rapid iteration, testing, and refinement of vehicle performance and safety features.
- **Cloud Service Providers (CSPs):** VAST provides scalable storage solutions for cloud providers, supporting multi-tenant environments with consistent performance. The DASE architecture enhances workload distribution and reliability, making it ideal for service providers managing diverse applications and customer requirements.
- **Energy and Utilities:** VAST's platform aids in seismic data analysis and smart grid monitoring, providing real-time data access for energy exploration and utility management. By supporting efficient resource monitoring and proactive infrastructure maintenance, VAST Data contributes to optimized resource allocation and cost savings in energy management.

- **Government:** Government agencies rely on VAST for mission-critical applications like disaster response modeling and defense research. Its high availability and real-time data processing capabilities enable secure and efficient data access, supporting critical decision-making in high-stakes environments.
- **Telecommunications:** For telecom providers, VAST enhances network optimization and customer experience analysis through its scalable, low-latency storage solution. The platform's architecture supports real-time analytics for network performance, allowing providers to respond rapidly and effectively to service demands.

VAST Data's industry-wide impact demonstrates its adaptability to diverse and demanding HPC applications, providing a unified, high-performance solution that drives productivity, enables cross-industry innovation, and supports future growth.

Simplified Administration and High Availability

One of VAST Data's strongest advantages is its ease of administration compared to conventional HPC storage solutions. With VAST Data, administrators can manage complex workflows without extensive tuning or third-party tools, thanks to features such as:

- **Centralized Management:** A unified interface provides centralized control over data placement, caching, load balancing, and performance monitoring, reducing the need for manual configurations.
- **Integrated High Availability and Snapshots:** The VAST Data Platform includes built-in high availability, automated failover, and snapshot capabilities, eliminating the need for additional tools to protect data and ensure uptime.
- **Reduced Complexity in Scaling:** VAST Data's architecture allows administrators to scale storage independently of computing, providing flexibility as data needs grow. Adding capacity doesn't require reconfiguring or tuning, streamlining scaling for HPC environments.

These features enable organizations to manage storage infrastructure more efficiently, reducing operational complexity and allowing administrators to focus on high-value tasks.

Future-Ready Storage for AI-Driven HPC

As HPC and AI continue to converge, organizations need infrastructure that can handle real-time data processing, vector embeddings, and immediate data access. The VAST Data Platform, particularly through its InsightEngine with NVIDIA, provides the foundation for advanced AI workflows that depend on large-scale data retrieval and processing.

- **Real-Time Data Access:** InsightEngine with NVIDIA processes data in real-time, embedding semantic relationships and enabling similarity searches across vast datasets. This capability supports advanced AI applications, accelerating decision-making and insights.

- **Scalable Semantic Database:** With support for trillions of vector embeddings, the VAST Data Platform can handle the extensive storage and processing needs of AI-driven HPC workloads, helping organizations leverage the full potential of their data.

By aligning with the demands of modern HPC and AI applications, VAST Data enables organizations to keep pace with technological advancements and prepare for future growth in data-driven research and development.

The Power of Unified HPC and AI Workloads

The simultaneous enablement of HPC and AI is exemplified in fields like quantitative trading, where ultra-low-latency execution and dynamic AI model training coexist. VAST Data's architecture empowers real-time execution of HPC workloads such as Monte Carlo simulations while seamlessly supporting AI-driven predictive analytics and sentiment analysis. By unifying HPC and AI on a single platform, VAST empowers industries to address hybrid workload demands without compromising performance, scalability, or simplicity.

- **Hybrid Workload Optimization:** VAST Data accelerates both HPC and AI workflows, ensuring seamless collaboration between low-latency HPC simulations and AI-driven predictive models.
- **Simplified Workflow Integration:** By eliminating data silos and supporting multiprotocol access, VAST enables industries to unify structured and unstructured data processing within a single architecture.

Key Takeaways for VAST Data's Value in HPC

To summarize, the VAST Data Platform offers the following benefits, making it an outstanding choice for HPC environments:

- **Unmatched Scalability:** VAST Data's DASE architecture and all-flash design allow organizations to scale seamlessly, supporting rapid data growth without the complexity of conventional storage systems.
- **High Performance and Low Latency:** With NVMe-based storage and support for RDMA and NVMe-oF, VAST provides the fast data access that is necessary for data-intensive HPC and AI applications.
- **Ease of Management:** VAST Data reduces administrative burden with automated data placement, multiprotocol support, and built-in data protection, eliminating many of the tuning and maintenance tasks associated with Lustre, and GPFS.
- **Future-Ready for AI Integration:** InsightEngine with NVIDIA provides real-time data processing and AI readiness, enabling organizations to confidently pursue advanced machine learning, AI, and data science initiatives.

By providing a platform that simplifies HPC data management, enhances performance, and prepares organizations for the future of AI and data-driven research, VAST Data addresses the core needs of modern HPC environments.

Conclusion

The VAST Data Platform represents a transformative leap in HPC storage capabilities, going beyond conventional solutions to address the core challenges modern HPC environments face. Its innovative architecture provides more than just speed and capacity, it also offers an adaptive, resilient platform that actively supports HPC and AI workloads. This enables organizations to achieve unprecedented levels of efficiency, flexibility, and insight. With the VAST Data Platform, users benefit from simplified administration, seamless scalability, and streamlined data access, allowing for more productive and cost-effective operations.

By integrating high-performance storage with real-time AI processing capabilities, VAST Data empowers researchers, engineers, and analysts to confidently tackle the most complex, data-intensive tasks. From scientific breakthroughs in genomics and climate modeling to advancements in engineering simulations and AI-driven analysis, VAST Data provides the robust infrastructure that high-impact research and development demand. The platform's multiprotocol support and DASE architecture removes many of the constraints seen with conventional HPC storage solutions like Lustre and GPFS. This gives organizations the agility to scale and innovate without the limitations of conventional architectures.

This white paper has explored how VAST Data supports HPC innovation by reducing the complexity of conventional data management and enhancing performance, reliability, and real-time processing. With its global namespace and integration with NVIDIA BlueField DPUs, VAST Data is not only equipped to meet today's data demands but also ready for the evolution and future convergence of HPC and AI. Features such as automated data placement, high availability, and real-time AI processing through InsightEngine with NVIDIA position VAST Data as a foundational solution that redefines what's possible in high-performance and data-driven computing.

As data volumes continue to expand and computational workloads grow in sophistication, the need for a next-generation data platform like VAST Data will only become more critical. With VAST Data, organizations can confidently face tomorrow's challenges, knowing they have the infrastructure to support both current demands and future aspirations. By choosing VAST, HPC environments are not just investing in advanced storage technology, they are embracing a platform built to power the discoveries, innovations, and transformations that define our digital future.



Are you ready to accelerate your HPC transformation?
Contact VAST Data today at hello@vastdata.com to learn
how our innovative platform simplifies data management,
maximizes performance, and future-proofs your most
demanding workloads.