

Architecting the Modern Data Platform: A Comparative Analysis of File, Object, and Table Storage for High-Scale Analytics

Ilya Malinovsky, Quanterra LLC

June 30, 2025

Executive Summary

The landscape of data infrastructure for high-scale analytics, particularly within computationally intensive domains like financial services, is undergoing a profound transformation. The exponential growth in data volume, velocity, and variety, coupled with the demands of artificial intelligence (AI) and machine learning (ML), has rendered traditional, monolithic data architectures inadequate. A new paradigm has emerged, centered on a layered, composite architecture where the choice of technology at each level—the physical storage platform, the on-disk data format, and the logical table format—is interdependent and must be driven by a nuanced understanding of specific workload characteristics. This report provides an analysis of this modern data stack, designed to guide technical leaders and architects in making strategic infrastructure decisions.

The analysis begins by deconstructing the three storage paradigms: file, object, and block storage. It explores their core architectures, access methods, and the fundamental trade-offs they present in terms of performance, scalability, and consistency. A central theme is the architectural tension between the POSIX API, the standard for traditional high-performance computing (HPC) and legacy applications, and the S3 API, the native language of scalable cloud storage. This conflict is a primary driver of innovation across the entire data stack.

Building on this foundation, the paper examines the critical role of analytical data formats. It details the superiority of columnar formats like Apache Parquet and ORC for optimizing data-at-rest through efficient compression and column pruning, which dramatically accelerates analytical queries. It then investigates Apache Arrow's role in optimizing data-in-motion, providing a standardized, zero-copy in-memory format that eliminates serialization bottlenecks between processing engines and programming languages.

The report then performs a deep dive into Apache Iceberg, the open table format that is revolutionizing the data lake. By imposing a transactional metadata layer on top of raw data files, Iceberg brings the reliability of a traditional

data warehouse-including ACID transactions, safe schema evolution, and point-in-time queries (time travel)-to the scalable, low-cost environment of the data lake, creating the foundation for the modern "Lakehouse" architecture.

Finally, the report provides a comparative, albeit non-exhaustive analysis of three leading high-performance data platforms that form the physical layer of this stack: the established and robust IBM Spectrum Scale (GPFS), the modern, software-defined Weka platform, and the revolutionary Disaggregated and Shared-Everything (DASE) architecture from VAST Data. The analysis reveals an architectural evolution from coupled systems to fully disaggregated platforms, a shift that offers unprecedented flexibility, scalability, and cost-efficiency. By synthesizing these layers, the report concludes with a strategic framework for designing a cohesive, high-performance analytics platform tailored to the demanding and diverse workloads of the modern data-driven enterprise.

1 Foundational Data Storage Paradigms

The selection of a data storage paradigm is the most fundamental architectural decision in designing a data platform. The choice between file, object, and block storage dictates not only how data is organized and accessed but also the inherent scalability, performance characteristics, and cost profile of the entire system. Understanding the architectural principles and trade-offs of each is essential before considering higher-level abstractions like data formats and table platforms.

1.1 File-Based Storage: The Hierarchical Model

File-based storage is the most familiar data organization paradigm, structuring data in a hierarchical system of directories and subdirectories that mirrors the logical organization of files on a personal computer. This intuitive model has been the bedrock of computing for decades, and its principles are deeply embedded in countless applications and operating systems.

Architecture and Access: Data is stored as discrete files, each located via a specific path within a directory tree (e.g., `/data/2024/market_data.csv`). This hierarchical structure is managed by a file system, which uses an index, often composed of inodes, to track the metadata and physical location of each file and directory. Access is typically provided through standard network protocols like Network File System (NFS) or Server Message Block (SMB), which are built upon the Portable Operating System Interface (POSIX) standard. POSIX defines a set of standard API calls for file operations (open, read, write, close, etc.), ensuring that applications written for one compliant operating system can function on another with minimal modification, a critical feature for software portability.

Metadata and Performance: In a traditional file system, metadata is relatively limited and standardized, typically including the file name, size, owner, permissions, and timestamps. While this is sufficient for general-purpose use,

the centralized nature of metadata management in traditional Network Attached Storage (NAS) systems creates a significant performance bottleneck at scale. As the number of files and directories grows into the millions or billions, traversing the hierarchical tree and managing the inode table becomes computationally expensive, leading to degraded performance and increased latency.

Evolution to Parallel File Systems (PFS): To address the scalability limitations of traditional NAS, Parallel File Systems (PFS) were developed for the HPC domain. Systems like IBM Spectrum Scale (formerly GPFS), BeeGFS (formerly FhGFS) and Lustre fundamentally re-architected the file system by distributing both data and metadata across multiple storage servers. Data for a single large file is "striped" across many servers and disks, allowing a client to read or write different parts of the file in parallel, achieving massive aggregate throughput. Similarly, metadata operations are distributed, eliminating the single-node bottleneck and enabling high levels of concurrent access from thousands of clients without significant performance degradation. This architecture makes PFS the standard for data-intensive HPC workloads in fields like financial modeling, scientific simulation, and media rendering.

1.2 Object-Based Storage: The Infinitely Scalable Model

Object-based storage represents a radical departure from the hierarchical file model, designed from the ground up for massive, exabyte-scale data repositories. It is the dominant storage paradigm in public cloud environments and the foundation for modern data lakes.

Architecture and Access: In an object store, data is managed as self-contained units called "objects" within a completely flat address space, often referred to as a storage pool or bucket. There is no concept of a directory hierarchy; each object is assigned a globally unique identifier that allows for its direct retrieval, regardless of how many other objects are in the system. This flat namespace is the key to its near-infinite horizontal scalability, as new storage nodes can be added to the pool without the need to rebalance a complex directory tree. Access to object storage is primarily through a RESTful HTTP-based API, with Amazon's Simple Storage Service (S3) API having become the de facto industry standard. This API-centric approach makes object storage natively accessible to modern cloud-native and web applications but requires legacy applications, which expect a POSIX file system interface, to be rewritten or to use a translation layer.

Metadata and Immutability: The most powerful feature of object storage is its support for rich, extensive, and customizable metadata. Unlike the fixed metadata of a file system, an object can be tagged with any amount of user-defined key-value pairs, which are stored alongside the data itself. This capability transforms the storage system from a passive repository into an active, queryable database of data assets, enabling sophisticated data management, indexing, and policy enforcement directly at the storage layer. For example, an X-ray image stored as an object could be tagged with patient ID, date, body part, and diagnosis, allowing for powerful, metadata-driven searches.

A core principle of object storage is the immutability of objects. An object cannot be partially modified in-place. To change an object, a new version of the entire object must be written. This atomic, version-based approach simplifies data protection and consistency in large distributed systems and is a foundational concept that enables the reliability features of higher-level platforms like Apache Iceberg.

1.3 Block-Based Storage: The Performance Foundation

Block-based storage is the lowest-level and highest-performance paradigm, serving as the fundamental building block for most other storage systems.

Architecture and Access: Data is broken down into fixed-size chunks called blocks, each with a unique address. The storage system itself is unaware of the concept of files or objects; it only manages these raw volumes of blocks. An operating system or a database management system layers a file system (like ext4 or NTFS) or a data structure on top of the block device to provide the logical organization that applications interact with. Metadata at the block level is minimal, limited to what is necessary for the system to locate and retrieve blocks, ensuring minimal overhead. Access is typically provided to a single host server over a high-speed Storage Area Network (SAN) using protocols like Fibre Channel or iSCSI. This direct, low-level access allows the host operating system to control the storage with maximum efficiency.

Performance: Block storage is engineered for one primary purpose: performance. It delivers extremely low latency and very high Input/Output Operations Per Second (IOPS). This makes it the ideal choice for performance-critical, transactional workloads such as job "scratch" space, relational databases, virtual machine disk images (VMDKs), and the underlying physical storage for high-performance file systems. While end-users and analytics applications rarely interact with block storage directly, its performance characteristics are what enable the responsiveness of the file and database systems built upon it.

1.4 The Great Divide: POSIX vs. S3 API

The divergence between the file and object storage paradigms is most clearly expressed in the deep architectural and semantic differences between their respective access protocols: POSIX and the S3 API. This is not merely a technical distinction; the tension between the requirements of the vast ecosystem of POSIX-dependent applications and the scalability advantages of the S3 object model is the central challenge that has driven the innovation of modern data platforms. Products like Apache Iceberg, Weka, and VAST Data are, in large part, sophisticated architectural responses to this fundamental conflict. They must either create a highly performant and scalable implementation of the POSIX standard or build new abstractions on top of the object model to provide the reliability guarantees that POSIX applications expect. More frequently than not, a modern product will enable access to the same data via S3 and POSIX, while implementing unified access controls.

Consistency Models: A primary point of divergence is the consistency model. POSIX mandates strong consistency, a strict guarantee that any modification to a file or directory is atomic and immediately visible to all processes and nodes in the cluster. This is essential for traditional applications, particularly databases and clustered applications that rely on shared file locking for coordination. However, enforcing strong consistency across a large, geographically distributed system is technically complex and can create performance bottlenecks. The S3 API, designed for web-scale distributed systems, originally offered an eventual-consistency model, where updates might take time to propagate across all replicas. While modern cloud object stores now typically provide strong read-after-write consistency for new objects, the guarantees for overwrites and deletes can still be weaker than what POSIX requires, a trade-off made to achieve massive scalability and availability.

Statefulness and Operational Semantics: POSIX is fundamentally a stateful protocol. The operating system kernel maintains state about open file descriptors, file pointers (the current position in a file), and file locks. This statefulness is extremely difficult to manage in a distributed, fault-tolerant manner. In contrast, the S3 API is designed to be stateless. Each HTTP request (e.g., GET, PUT, DELETE) is a self-contained, atomic operation that includes all necessary information, simplifying load balancing and failure recovery in large clusters. This leads to significant differences in supported operations. POSIX offers a rich set of semantics, including the ability to modify files in-place (e.g., appending data or overwriting a specific byte range) and sophisticated file locking mechanisms. These capabilities are assumed by countless scientific, financial, and engineering applications. The S3 API is far more restrictive, supporting only whole-object operations. This immutability model simplifies the backend implementation but forces a complete rewrite of applications that rely on in-place file modifications.

Bridging the Gap: The desire to run POSIX-compliant applications on cost-effective and scalable object storage has led to the creation of translation layers like s3fs-fuse. These tools present an S3 bucket as a local file system. However, they often suffer from severe performance limitations and compatibility issues because they must emulate complex POSIX semantics (like directory listings and file appends) using inefficient sequences of S3 API calls. This fundamental impedance mismatch highlights that a truly high-performance solution cannot simply translate between the two worlds; it must be architected from the ground up to resolve these core tensions.

2 Data Formats for High-Throughput Analytics

While the storage paradigm defines the macro-level organization of data, the internal format of the data files themselves is equally critical for the performance of analytical workloads. The choice of file format dictates how efficiently data can be read from storage, compressed to save space, and processed in memory. The industry has seen a decisive shift away from traditional row-based formats

toward columnar formats for data-at-rest, and the emergence of a standardized in-memory columnar format to accelerate data-in-motion.

2.1 The Columnar Advantage: Apache Parquet and ORC

For decades, data was primarily stored in row-oriented formats like Comma-Separated Values (CSV) or Avro. In these formats, all the fields for a single record are stored contiguously on disk. This is efficient for transactional systems (OLTP), where the common operation is to read or write an entire record at once. However, for analytical queries (OLAP), this model is profoundly inefficient. An analytical query, such as calculating the average trade size, moving average or building a book for a specific stock, typically only needs to access a few columns (e.g., `ticker`, `trade_size`) from a table with potentially hundreds of columns. In a row-based format, the query engine must read every column for every row from disk, discarding the vast majority of the data it just read, leading to massive I/O inefficiency.

Columnar storage formats, most notably Apache Parquet and Apache ORC, were created to solve this specific problem. They store data by column rather than by row, meaning all values for the `ticker` column are stored together, all values for the `trade_size` column are stored together, and so on. This simple change in layout provides two transformative benefits for analytics.

Performance Gains from Column Pruning: The primary advantage of columnar formats is column pruning (also known as predicate pushdown). Predicate pushdown is a query optimization technique that significantly improves query performance by moving filtering conditions (predicates) closer to the data source. Instead of loading the entire dataset into memory and then filtering it, the database pushes down the filtering logic to the data source (like a database table or storage layer), allowing the source to filter the data before sending it back. This reduces the amount of data that needs to be transferred and processed, leading to faster query execution and more efficient resource utilization.

Compression Efficiency: Data within a single column is homogeneous - it is all of the same data type and often shares similar characteristics and value ranges. This uniformity allows for much more effective compression than is possible with row-based formats, where disparate data types are mixed together. Both Parquet and ORC leverage modern compression algorithms like Snappy and Zstandard (zstd) to reduce storage footprints by up to 75 or more, which translates directly into lower storage costs, IO and bandwidth consumption.

Parquet vs. ORC Comparison: While both Parquet and ORC are powerful columnar formats, they have different origins and subtle architectural trade-offs.

- **Apache Parquet:** Developed by Cloudera and Twitter, Parquet has become the de facto standard columnar format across the big data ecosystem, with strong support in Apache Spark, AWS Athena, Google BigQuery, and many other platforms. It excels at handling complex, nested data struc-

tures (e.g., JSON-like data within a column) and provides robust support for schema evolution, allowing columns to be added, removed, or renamed over time without requiring a complete rewrite of the dataset. This flexibility is crucial in dynamic analytical environments where data sources and requirements frequently change.

- **Apache ORC (Optimized Row Columnar):** Originally developed within the Apache Hive project, ORC was designed for optimal performance in Hadoop-based data warehouses. ORC often achieves slightly higher compression ratios than Parquet and can deliver faster query performance for aggregation-heavy workloads. This is due to its inclusion of lightweight indexes, or "stripe-level statistics" (min/max values, counts), within the file itself. These indexes allow the query engine to skip entire blocks of data within a column, further reducing I/O. ORC also supports vectorized processing, which allows operations to be performed on batches of rows at a time, improving CPU efficiency. However, its schema evolution capabilities are generally considered less flexible than Parquet's.

While these formats are transformative for most analytical workloads, they are not a panacea. For extremely wide and sparse tables, such as those found in some machine learning feature stores, the metadata required to describe hundreds or thousands of columns can become large and complex. In these edge cases, the time spent reading and parsing the file metadata can begin to dominate the overall query time, creating a new bottleneck.

2.2 Accelerating In-Memory Analytics: Apache Arrow

While Parquet and ORC solve the problem of efficiently reading data from persistent storage (data-at-rest), a significant performance bottleneck remained in the processing of that data once it was loaded into memory (data-in-motion). Historically, each data processing system (like Spark), programming language (like Python), and library (like Pandas) had its own proprietary in-memory data representation.

This created what can be called the "last mile" problem of data analytics. To move data from a Spark cluster to a Python process for machine learning with Pandas or NumPy, the data had to be serialized from Spark's internal format into an intermediate format (like JSON or pickle), sent over the network, and then deserialized into Python's internal format. This serialization and deserialization (often shortened to SerDe) process is computationally expensive, consuming significant CPU cycles and memory bandwidth, and often becoming the primary bottleneck in data-intensive pipelines.

Arrow's Solution: A Standardized In-Memory Format: Apache Arrow was created to solve this problem by defining a standardized, language-independent columnar in-memory format. It is not a file format for long-term storage like Parquet; rather, it is a specification for how structured, tabular data should be laid out in RAM for optimal processing on modern hardware.

Zero-Copy Data Sharing and Hardware Optimization: The core benefit of Arrow is enabling zero-copy data sharing. When two different systems or processes both "speak" Arrow, they can exchange data by simply passing a memory pointer. There is no need to copy, serialize, or deserialize the data, as both processes can read and operate on the exact same block of memory. This eliminates the SerDes overhead and can accelerate data interchange performance by 100x or more. Furthermore, the Arrow memory layout is specifically designed to be friendly to modern CPUs and GPUs. Its contiguous columnar layout improves CPU cache locality and enables the use of SIMD (Single Instruction, Multiple Data) instructions, which perform the same operation on multiple data points simultaneously, leading to highly efficient parallel processing.

The Arrow Ecosystem: The Arrow project extends beyond just the memory format to include a suite of tools for building high-performance data applications:

- **Arrow Flight:** A high-performance, parallel-aware RPC (Remote Procedure Call) framework built on gRPC. It is designed for bulk transfer of Arrow data between systems over a network, avoiding the bottlenecks of traditional protocols like ODBC/JDBC by eliminating SerDe overhead.
- **Plasma Object Store:** A shared-memory object store that allows multiple, independent processes on the same machine to access and share large Arrow data structures with zero copying. A data-loading process can place an Arrow table in Plasma, and multiple data science or analytics processes can then access that data instantly without incurring any copy or deserialization costs.

The adoption of Arrow creates a powerful synergy with formats like Parquet. A modern analytics pipeline can be designed to efficiently read specific columns from a Parquet file on disk directly into the Arrow in-memory format. This Arrow data can then be passed frictionlessly between different processing stages, such as a Spark SQL engine and a Python-based machine learning library, achieving end-to-end performance optimization.

3 The Data Lakehouse Revolution: Apache Iceberg

The rise of scalable object storage and powerful analytical formats like Parquet enabled the creation of massive data lakes. However, these data lakes often devolved into unreliable "data swamps" because they lacked the critical data management features of traditional data warehouses, such as transactional integrity, schema enforcement, and reliable performance. Apache Iceberg is an open-source table format designed to solve this problem by bringing the reliability and simplicity of SQL tables to the vast, low-cost datasets stored in data lakes, thereby enabling the "Lakehouse" architecture.

3.1 Iceberg Architecture: A Metadata-Driven Approach

The fundamental flaw of early data lake table formats (like those used by Apache Hive) was that the table was defined simply as a collection of directories and files. The query engine had to discover the state of the table by performing expensive and slow file system listing operations, and any operation that was not atomic could leave the table in a corrupt, inconsistent state.

Iceberg’s innovation is to decouple the logical table from the physical data files by introducing a formal, multi-level metadata layer. This layer provides a complete and consistent definition of the table’s state at any point in time, and all operations are performed by manipulating this metadata atomically. The architecture consists of a hierarchy of metadata files that provide a progressively more detailed view of the table’s contents.

- **The Catalog:** This is the entry point for any query engine. The catalog is a simple key-value store that maps a table name (e.g., `prod.market_data.trades`) to the location of the single, current top-level metadata file for that table. This can be an existing service like the Hive Metastore or AWS Glue, or a database via JDBC. The atomic “swap” of this pointer is the mechanism that enables ACID transactions.
- **Metadata Files (Snapshots):** The file pointed to by the catalog is the current metadata file. Each metadata file represents an immutable snapshot of the table at a specific point in time. It contains the table’s schema at that time, its partitioning configuration, and a pointer to a manifest list.
- **Manifest Lists:** A manifest list contains pointers to one or more manifest files. Crucially, it also stores summary metadata about each manifest file, such as the range of partition values contained within it. This allows a query engine to quickly prune entire manifest files (and thus all the data files they point to) that are not relevant to a query’s filter conditions.
- **Manifest Files:** Each manifest file contains a list of the actual data files (e.g., Parquet or ORC files) that comprise the table’s data. For each data file, the manifest stores its path, its partition data, and detailed column-level statistics, such as the minimum and maximum values for each column within that file. These statistics are the key to Iceberg’s performance, as they enable the query engine to perform fine-grained file pruning.

All changes to an Iceberg table—whether inserting data, deleting rows, or altering the schema—are executed through a consistent, atomic commit process: new data and metadata files are written, and then the catalog’s pointer is atomically updated to point to the new top-level metadata file. The old files are not immediately deleted, which is what enables features like time travel and rollback.

3.2 Bringing Reliability to the Data Lake

The metadata-driven architecture of Iceberg enables a suite of powerful features that were previously exclusive to traditional data warehouses. This reliability is often best understood through the analogy of a version control system like Git. Each commit to an Iceberg table is like a git commit—it creates a new, immutable snapshot of the data’s state. The catalog acts like a branch pointer (e.g., `main`), always pointing to the latest version. This model provides the foundation for DataOps and MLOps, enabling versioning, reproducibility, and safe experimentation on production data.

- **ACID Transactions:** Iceberg provides full ACID (Atomicity, Consistency, Isolation, Durability) guarantees. Every operation is atomic; it either completely succeeds or completely fails, leaving the table in a consistent state. This prevents data corruption from partially failed write jobs, a common problem in traditional data lakes. For financial institutions, this transactional integrity is non-negotiable for managing critical datasets like trade records or customer accounts.
- **Full Schema Evolution:** Iceberg supports a full range of schema evolution operations—adding, dropping, renaming, reordering, or changing the type of columns—without requiring a costly and disruptive rewrite of the underlying data files. Iceberg tracks schema changes in its metadata, allowing it to correctly interpret data files written with older schemas. This agility is vital in financial analytics, where new data sources or regulatory requirements can necessitate frequent changes to data structures.
- **Hidden Partitioning and Partition Evolution:** Iceberg abstracts the physical data layout from the logical partitioning scheme. Users can define partitions based on transformations of a column (e.g., partitioning a trade table by `month(trade.timestamp)`), and Iceberg handles the value generation automatically. This is called hidden partitioning because users query the raw timestamp column without needing to know the partition structure. More powerfully, the partitioning scheme itself can be changed over time (partition evolution) without rewriting old data. For example, a table initially partitioned by month can be changed to be partitioned by day for new data, and Iceberg will seamlessly handle queries that span both partitioning schemes.
- **Time Travel and Rollback:** Because Iceberg never modifies data in-place and retains a history of metadata snapshots, users can query the table as it existed at any previous snapshot ID or timestamp. This time travel capability is invaluable for reproducible research, allowing analysts to rerun a model on the exact same version of the data that was used previously. It is also a powerful tool for auditing and compliance. Furthermore, if a bad data load corrupts the table, an administrator can perform an instantaneous rollback by simply resetting the catalog’s pointer to a

previous, valid snapshot, providing a powerful and fast disaster recovery mechanism.

3.3 Performance and Ecosystem Integration

Beyond reliability, Iceberg is engineered for high performance at massive scale, primarily by providing rich metadata to query engines so they can minimize the amount of data they need to read.

Query Performance Optimization: Iceberg’s primary performance benefit comes from data pruning. Before a query runs, the engine consults Iceberg’s metadata. First, it uses the partition information in the manifest lists to eliminate any manifest files whose partition value ranges do not overlap with the query’s filters. Then, for the remaining manifests, it uses the column-level min/max statistics to eliminate any individual data files that cannot possibly contain matching rows. This multi-level pruning can dramatically reduce the data scanned, with benchmarks showing query performance improvements of up to 50 percent compared to querying raw Parquet files.

Eliminating File System Bottlenecks: By maintaining a complete list of data files in its manifest files, Iceberg eliminates the need for query engines to perform slow and expensive recursive directory listings on the underlying storage system. This is a particularly significant advantage on cloud object stores like Amazon S3, where LIST operations can be slow and are subject to API rate limiting, which often causes job failures at scale.

Engine Integration: A core design principle of Iceberg is to be an open, engine-agnostic standard. It provides well-defined Java and Python libraries for interacting with its metadata, allowing a wide range of processing engines—including Apache Spark, Trino, Flink, Presto, and Hive—to read from and write to the same Iceberg tables concurrently and safely. This interoperability is a key advantage, as it allows an organization to build a single, unified data lakehouse and empower different teams to use the best tool for their specific job—Spark for large-scale ETL and ML, and Trino for interactive SQL analytics—all on the same source of truth.

4 A Comparative Analysis of High-Performance Data Platforms

While table formats like Iceberg and data formats like Parquet define the logical structure and encoding of data, the performance, scalability, and reliability of an analytics platform ultimately depend on the underlying physical storage system. This section provides an architectural comparison of three leading platforms in the high-performance data space: the long-standing HPC champion, IBM Spectrum Scale (GPFS); the modern, software-defined Weka platform; and the revolutionary disaggregated architecture of VAST Data. These systems represent three distinct generations of storage architecture, each with unique design

philosophies that make them suitable for different types of large-scale financial and AI workloads.

4.1 IBM Spectrum Scale (GPFS): The Legacy of Parallelism

IBM Spectrum Scale, originally known as the General Parallel File System (GPFS), is a mature and field-proven parallel file system that has been the backbone of many of the world's largest supercomputers and enterprise HPC environments for decades.

Architecture: GPFS is based on a shared-disk clustered architecture. In this model, all nodes in the cluster have direct, concurrent block-level access to a common set of storage devices over a high-speed fabric like Fibre Channel SAN or, more recently, fast Ethernet. The file system logic runs on all nodes, and a sophisticated distributed lock and token management system is used to maintain data coherency and consistency, allowing thousands of clients to read and write to the same files simultaneously without corruption. Data is striped across numerous disks and storage servers to enable parallel I/O, delivering the extreme throughput required for traditional HPC applications.

Protocol Support and Features: At its core, Spectrum Scale is a POSIX-compliant file system, making it an ideal platform for the vast ecosystem of scientific and engineering applications that rely on this standard. Over time, IBM has extended its capabilities to serve a broader range of enterprise needs. Through its Cluster Export Services (CES) nodes, it can also serve data via standard NAS protocols (NFS, SMB) and object storage protocols (S3, Swift), providing a unified storage solution. For financial analytics, several features are particularly relevant:

- **Global Namespace and Active File Management (AFM):** Spectrum Scale can create a single global namespace that spans multiple physical sites and even cloud object storage. AFM allows a local cluster to act as a high-performance cache for data residing in a remote cluster or an S3 bucket. This enables hybrid cloud workflows, where on-premises compute can "burst" to the cloud for peak demand, and facilitates multi-site collaboration and disaster recovery-critical capabilities for global financial institutions.
- **Enterprise-Grade Reliability:** Having been deployed in the most demanding environments for over two decades, Spectrum Scale is known for its robustness, data integrity features, high availability, and sophisticated data management tools, including policy-driven tiering and snapshots. Several Fortune 500 financial service providers rely on it for mission-critical applications.

Modernization and Limitations: Spectrum Scale has evolved to integrate with modern data ecosystems, offering a Hadoop connector for in-place

analytics with platforms like Cloudera and containerized storage drivers for Kubernetes environments. However, its architecture, originally conceived in the era of spinning disks (HDDs), can be complex to deploy and manage, often requiring deep specialized expertise. While it performs well on flash, it was not designed as a "flash-native" system, and its coupled architecture can be less flexible and cost-efficient than newer, software-defined approaches.

4.2 Weka: A Software-Defined Architecture for AI

Weka has emerged as a leader in the new generation of parallel file systems, with an architecture designed from the ground up for the performance characteristics of NVMe flash storage and the mixed, demanding I/O patterns of AI and machine learning workloads.

Architecture: WekaFS is a fully distributed, software-defined parallel file system built on a unique container-based architecture. The core Weka software processes run in user space within lightweight Linux Containers (LXC), which are assigned dedicated CPU cores and memory. This approach bypasses the Linux kernel for most I/O operations, minimizing context switching and system call overhead, which results in extremely low and predictable latency. The architecture is fully distributed, using hashing algorithms to spread both data and metadata evenly across all nodes in the cluster, which eliminates the hotspots and bottlenecks common in legacy NAS systems.

Protocol Support and Features: Weka provides a multi-protocol data platform with a single, unified namespace. It offers a high-performance, fully POSIX-compliant client for demanding HPC and AI applications, alongside standard NFS, SMB, and S3 interfaces for broader compatibility. This allows a single platform to serve GPU-based training clusters, data ingest pipelines, and analytics users simultaneously. Key features for finance and AI include:

- **Flash-Native Performance:** WekaFS is specifically optimized for the low latency and high parallelism of NVMe SSDs. It is capable of delivering tens of gigabytes per second of throughput to a single client, saturating the network links to powerful GPU servers and dramatically reducing I/O wait times during ML model training.
- **Integrated, Transparent Tiering:** A core feature of the Weka architecture is its ability to seamlessly and automatically tier data from its high-performance flash tier to any S3-compatible object store, whether on-premises or in the public cloud. This tiering is completely transparent to applications; the entire dataset resides in a single namespace. This allows organizations to build a relatively small, high-performance flash tier for active data while cost-effectively scaling capacity into the petabytes on a cheaper object storage backend.
- **Cloud-Native and Hybrid Flexibility:** The Weka software is hardware-agnostic and can run on standard x86 servers, both on-premises and in any major public cloud. This provides maximum flexibility, enabling easy

migration of workloads, cloud bursting for peak demand, and consistent data management across hybrid environments.

Limitations: As a high-performance system, Weka achieves its best results with high-speed, low-latency networking like InfiniBand with RDMA, which can add cost and complexity to the infrastructure. It is also a proprietary, licensed software product.

4.3 VAST Data: The Disaggregated, Shared-Everything (DASE) Paradigm

VAST Data introduced a fundamentally new storage architecture that breaks with the decades-old principles of shared-nothing and shared-disk systems. This Disaggregated and Shared-Everything (DASE) architecture was designed to eliminate the trade-offs between performance, scalability, and cost that have plagued legacy systems.

Architecture: The DASE paradigm is built on two revolutionary concepts:

1. **Disaggregation:** The compute logic of the storage system is physically separated from the storage media. The system is composed of two types of nodes connected by a high-speed NVMe-over-Fabrics (NVMe-oF) network: stateless C-nodes (Compute Nodes), which are servers running the VAST software in containers and handling all protocol requests (NFS, SMB, S3), and stateful D-nodes (Data Nodes), which are simple, highly-available NVMe enclosures that hold the storage media.
2. **Shared Everything:** Every C-node has direct, low-latency access to every NVMe drive in every D-node across the NVMe-oF fabric. This is a crucial distinction from shared-nothing architectures, where a node owns its local drives. In the DASE model, there is no need for performance-killing "east-west" traffic between compute nodes to coordinate I/O. All metadata is shared and accessible to all C-nodes, allowing for massively parallel and lock-free access to the underlying storage. This eliminates the bottlenecks that cause performance to degrade as traditional scale-out clusters grow.

This architecture allows for the independent scaling of compute and capacity. If more performance is needed (e.g., more NFS or S3 throughput), an organization can simply add more C-nodes. If more capacity is needed, they can add more D-nodes. This provides unprecedented flexibility and cost-efficiency, allowing infrastructure to be precisely tailored to workload requirements.

Storage Media and Data Reduction: VAST's architecture is designed to make all-flash storage affordable for all data. It achieves this by using low-cost, high-density QLC flash for bulk data storage, paired with a small amount of higher-endurance Storage Class Memory (SCM) to act as a persistent write buffer and to store metadata. This is combined with a groundbreaking similarity-based data reduction algorithm that finds and eliminates redundant

patterns globally across the entire dataset, achieving dramatically higher reduction rates than traditional compression and deduplication, and bringing the effective cost of flash down to levels competitive with hard disk-based archive systems.

Protocol Support and Platform Extension: VAST provides a unified data platform that serves file (NFS, SMB), object (S3), and block protocols from a single, tierless pool of storage. It has further extended this platform with the VAST DataBase, an ACID-compliant transactional and analytical database built directly into the storage system, and the VAST DataSpace, a global namespace that provides consistent, performant data access across geographically distributed clusters.

5 Synthesis and Strategic Recommendations

The preceding sections have deconstructed the modern data stack into its constituent layers: the foundational storage paradigms, the analytical data formats, the transactional table format, and the high-performance physical platforms. This final section synthesizes these elements, illustrating how they combine to form a cohesive architecture and providing a strategic framework for making informed decisions based on specific workload requirements. The analysis culminates in a forward-looking perspective on the future trajectory of data infrastructure.

5.1 Building the Modern Data Stack: A Layered Approach

The true power of the modern data stack lies not in any single component but in the synergistic integration of technologies at each layer. The architecture is characterized by a deliberate decoupling of concerns: the physical storage is decoupled from the compute engines, and the logical table structure is decoupled from the physical file layout. This separation allows for unprecedented flexibility, performance, and scalability.

To illustrate this, consider a common, yet demanding, workflow at a quantitative hedge fund: backtesting a new algorithmic trading strategy against a decade of high-resolution historical tick-by-tick data and subsequently training a machine learning model to refine its signals.

1. **Storage Platform Layer:** The raw data, consisting of petabytes of tick-by-tick market data, is stored on a high-performance data platform like VAST Data or Weka. This platform provides a single, scalable namespace accessible via both NFS for traditional research tools and a high-performance S3 API for modern data science frameworks. Its all-flash architecture ensures that any slice of the historical data can be accessed with low latency, eliminating the I/O bottlenecks that would plague a traditional tiered storage system.

2. **Table Format Layer:** This huge collection of raw Parquet files is organized and managed as an Apache Iceberg table. This layer is critical for several reasons. First, it provides ACID transactions, allowing data engineers to reliably correct historical data for corporate actions like stock splits or dividend payments without taking the system offline or risking data corruption. Second, its time travel capability is essential for point-in-time correctness; a backtest for a strategy in 2015 can be run against the exact snapshot of the data as it existed on that day, preventing lookahead bias. Third, as new data sources are added (e.g., alternative data), schema evolution allows the table to be updated without invalidating historical data. It also supports migrating the data across physical platforms without losing its context.
3. **Processing Engine Layer:** An Apache Spark cluster is used to execute the backtest. When the query is submitted, Spark’s planner interacts with the Iceberg catalog and metadata. Iceberg’s multi-level metadata allows Spark to perform aggressive data pruning, identifying and reading only the specific Parquet files and columns relevant to the tickers and time period of the backtest, dramatically reducing the amount of data that needs to be processed.
4. **In-Memory Format Layer:** As Spark processes the data, it uses Apache Arrow to represent the data in-memory. When the results of the backtest are passed to a Python-based machine learning framework (like TensorFlow or PyTorch) for model training, the data is transferred between the Java-based Spark executors and the Python processes in the Arrow format. This zero-copy transfer eliminates computationally expensive serialization and deserialization steps, ensuring that the high-performance GPU servers are fed with data as quickly as possible, maximizing their utilization.

This integrated pipeline demonstrates the power of the layered, decoupled approach. Each component is optimized for its specific task, and open standards (S3, Parquet, Iceberg, Arrow) ensure seamless interoperability, creating a whole that is far greater than the sum of its parts.

5.2 Architectural Decision Matrix

Selecting the right combination of technologies requires a clear understanding of the specific workload characteristics and business objectives. There is no single "best" architecture; the optimal choice is a function of the trade-offs an organization is willing to make. Architects should consider the following dimensions:

Workload Characterization:

- **Read-Heavy Analytics and Backtesting:** These workloads are characterized by large, sequential reads of historical data. Performance is primarily gated by storage throughput and the ability of the query engine to

prune data effectively. The combination of a columnar format (Parquet), a table format with rich statistics (Iceberg), and a high-throughput storage platform (any of the three discussed) is ideal. The key is to minimize I/O.

- **Write-Intensive Data Ingest:** This involves capturing real-time or near-real-time data streams, such as market data feeds. This pattern often generates many small files and requires low write latency. Here, the choice of data pipeline pattern-ETL (Extract, Transform, Load) vs. ELT (Extract, Load, Transform)-becomes critical. An ELT pattern, where raw data is landed quickly on the storage platform and transformed later, is often preferred. The storage system must be efficient at handling small file writes and subsequent compaction operations, a strength of modern log-structured file systems.
- **Mixed AI/ML Workloads:** This is often the most demanding workload, involving a mix of sequential reads (for training data), random reads (for feature lookups), and small file writes (for checkpoints and logs). It requires a platform that can deliver high performance across all I/O patterns and provide low-latency access to GPU servers. Modern, flash-native platforms like Weka and VAST are specifically architected for these mixed workloads.

Strategic Decision Factors:

- **Primary Access Protocol:** Does the application ecosystem rely heavily on POSIX-compliant file access, or is it built around the S3 API?. If POSIX is non-negotiable, a high-performance parallel file system like Spectrum Scale or Weka is a natural fit. If the architecture is cloud-native and S3-centric, a platform with a strong S3 implementation like VAST Data may be preferable.
- **Data Scale and Growth:** The expected scale of the data-terabytes, petabytes, or exabytes-will influence the choice. While all three platforms can scale to petabytes, architectures like VAST's DASE are designed with exabyte-scale in mind, where eliminating inter-node chatter becomes paramount for performance.
- **Cloud Strategy:** The organization's cloud strategy is a critical factor. Is the goal to remain on-premises, build a hybrid cloud with bursting capabilities, or migrate fully to the cloud?. Platforms like Weka and VAST, with software that runs consistently across all environments, offer the most flexibility for hybrid and multi-cloud strategies.
- **Total Cost of Ownership (TCO):** A comprehensive TCO analysis must go beyond the initial acquisition cost to include power, cooling, data center space, and operational overhead. The ability of platforms like VAST to use low-cost QLC flash and achieve high data reduction can significantly

impact long-term TCO. Similarly, the ability to scale compute and capacity independently offers a more granular, pay-for-what-you-need economic model compared to traditional hyperconverged systems.

5.3 Future Outlook: The Converged, Data-Defined Future

The trends and technologies analyzed in this report point toward a clear future for high-performance data infrastructure, defined by convergence, abstraction, and simplification.

The End of Storage Tiering: A significant implication of architectures like VAST Data's is the potential obsolescence of complex, multi-tiered storage. Historically, infrastructure was built with a hot tier of expensive, fast storage (NVMe), a warm tier of less expensive SSDs, and a cold archive tier of cheap HDDs or object storage. This created immense operational complexity, requiring data to be constantly moved between tiers, which introduced latency and created bottlenecks. By using low-cost QLC flash combined with revolutionary data reduction, VAST makes a single, tierless all-flash platform economically viable for all data, from active to archive. Weka achieves a similar outcome through its transparent tiering to an object store, presenting a single namespace to the user. This paradigm shift radically simplifies data architecture. It eliminates data movement and makes the entire historical dataset instantly accessible for analysis, providing a significant competitive advantage to financial firms that can now train models and run analytics on their most valuable asset—their complete data history—without delay.

The Rise of the Data Lakehouse: The combination of a scalable, high-performance physical storage platform with a transactional, open table format like Apache Iceberg is the definitive architecture for the future of data analytics. The Lakehouse paradigm successfully merges the low-cost scalability and flexibility of the data lake with the reliability, performance, and data management capabilities of the traditional data warehouse. This unified approach breaks down the silos between data engineering, data science, and business intelligence, allowing all teams to work on a single, consistent, and reliable source of data.

From Storage to Data Platform: The most advanced systems are evolving beyond being passive repositories for data. Platforms like VAST, with its integrated DataBase and DataEngine, are starting to embed query logic, event triggers, and computational functions directly into the storage layer. This trend toward a "data-defined" computing architecture, where processing moves closer to the data, will continue to blur the lines between storage and compute. The result will be more efficient, intelligent, and autonomous data platforms that not only store data but also actively participate in its processing and analysis, forming the dynamic foundation for the next generation of AI and data-intensive applications.

Table 1: Comparison of Foundational Storage Paradigms

Feature	File Storage	Object Storage	Block Storage
Architecture	Hierarchical tree of directories and files.	Flat address space of self-contained "objects" in a storage pool.	Raw volumes of fixed-size blocks with no inherent structure.
Primary Access Protocol	POSIX-based (NFS, SMB).	RESTful HTTP API (S3 API is the de facto standard).	Low-level storage protocols (Fibre Channel, iSCSI, NVMe-oF).
Metadata and Consistency Model	Limited, standardized metadata (name, size, permissions) managed by the file system's inode table. Strong consistency; changes are immediately visible to all clients.	Rich, extensive, and user-customizable key-value metadata stored with each object. Typically eventual consistency, though modern systems offer strong read-after-write consistency for new objects.	Minimal metadata, typically just the block address, to reduce overhead. Managed by the overlying file system or database; provides strong consistency.
Scalability and Performance Profile	Traditional NAS scales vertically and is limited. Parallel File Systems scale horizontally to petabytes. Low latency for small file I/O. PFS provides very high throughput for large files.	Near-infinite horizontal scalability to exabytes and trillions of objects. Higher latency due to HTTP overhead, but high throughput for large objects.	Scales by adding more volumes; performance can be a bottleneck at extreme scale. Lowest latency and highest IOPS. Optimized for transactional performance.
Key Operations and Ideal Workload	Create, read, write, delete, append, seek, lock. Supports in-place modification. Shared corporate documents, home directories, HPC simulations, applications requiring POSIX compliance.	PUT, GET, DELETE, LIST. Objects are immutable; modifications require writing a new version. Cloud-native applications, data lakes, archives, backups, rich media storage, large-scale unstructured data analytics.	read block, write block. Raw access provided to a host operating system. Transactional databases (SQL/NoSQL), virtual machine disks, performance-critical enterprise applications.

Table 2: Feature Comparison of Analytical Data Formats

Feature	Apache Parquet	Apache ORC	Apache Arrow
Primary Use Case	Efficient, long-term storage of analytical data at rest.	High-compression, long-term storage, especially in Hadoop/Hive ecosystems.	High-performance, in-memory data representation and interchange (data in motion).
Data Structure	Columnar, on-disk format. Supports complex and nested data types.	Columnar, on-disk format. Includes built-in lightweight indexes (stripe statistics).	Language-independent, columnar in-memory format. Not a persistent storage format.
Compression	High compression efficiency using algorithms like Snappy, Gzip, Zstd.	Very high compression efficiency, often slightly better than Parquet due to its structure.	Supports compression for in-memory data, but primary focus is on processing speed, not storage size.
Schema Evolution	Robust support for adding, removing, renaming, and re-ordering columns without data rewrite.	Supports schema evolution, but generally considered less flexible than Parquet.	Schema is part of the in-memory data structure; not designed for long-term evolution in the same way as storage formats.
Mutability	Immutable files. Changes require creating new files.	Immutable files. Changes require creating new files.	Immutable data buffers. Designed for read-heavy analytical processing.
Key Performance Feature	Column pruning (predicate push-down) to minimize I/O from disk.	Column pruning and stripe-level statistics to skip data blocks within files.	Zero-copy data sharing between processes and systems, eliminating serialization/deserialization overhead.
Ecosystem Integration	Very broad support across Spark, Presto, Trino, Dremio, AWS, Google Cloud, etc.	Strongest within the Hadoop ecosystem (Hive, Presto), but also supported by Spark and others.	Integrated into a vast and growing number of projects including Spark, Pandas, Dremio, and many databases for data interchange.

Table 3: Architectural and Feature Comparison of High-Performance Data Platforms

Feature	IBM Spectrum Scale (GPFS)	Weka	VAST Data
Core Architecture Scalability Model	Shared-disk, clustered parallel file system. Scales out by adding nodes to the cluster. Can scale performance and capacity together.	Software-defined, NVMe-native parallel file system with a containerized, user-space I/O path. Scales out by adding standard x86 servers. Performance scales linearly with the cluster size.	Disaggregated and Shared-Everything (DASE); stateless compute nodes (C-nodes) separated from stateful storage nodes (D-nodes) via NVMe-oF. Independent scaling of compute (C-nodes) and capacity (D-nodes).
Primary Media	Supports HDD, SSD, and Flash. Originally designed for HDD performance.	NVMe-native. Optimized exclusively for flash performance.	Low-cost QLC Flash for data, with Storage Class Memory (SCM) for write buffering and metadata.
Protocol Support	POSIX, NFS, SMB, S3, HDFS. Protocols are often served by dedicated gateway nodes (CES).	POSIX, NFS, SMB, S3. All protocols served natively by all nodes.	NFS, SMB, S3, Block, and a native Database API. All protocols served by all C-nodes from a single data store.
Data Reduction	Compression.	Standard compression and deduplication.	Global, similarity-based data reduction, delivering very high efficiency across all data types.
Resiliency Model	Replication and network-level RAID (declustered RAID in ESS appliances).	Distributed, grid-based erasure coding designed for flash; fast rebuilds.	Locally decodable, grid-based erasure codes written in a log-structured format to maximize QLC endurance.
Metadata Handling	Distributed lock and token management system.	Fully distributed metadata across all nodes; no single point of failure.	Stored on high-performance SCM and accessible to all C-nodes in a shared-everything model.
Cloud Integration	Hybrid cloud via Active File Management (AFM) for caching and tiering.	Native software runs on-prem or in any public cloud. Transparent, policy-based tiering to object storage.	Native software can run in the cloud. VAST DataSpace provides a global namespace across on-prem and cloud deployments.
Ideal Workloads	Traditional HPC, large-scale sequential I/O, enterprise file sharing environments	AI/ML training and inference, genomics, financial analytics, and mixed I/O workloads	At-scale data lakes and AI clouds, quantitative research on massive historical datasets,