

Want to Maximize NVMe Performance? Keep Your Eye on Operational Efficiency

The 451 Take

For datacenter storage, keeping storage and server compute resources separate is a prime virtue, if not a directive, because of its multiple benefits in terms of [operational flexibility and costs](#). This is called disaggregation, and it involves the use of stand-alone storage systems that are shared by multiple application servers, which is far more efficient than the use of isolated storage capacity within each application server. Also known as networked (SAN or NAS) storage, disaggregated storage has become the dominant form of storage in both large and small datacenters over the last two decades.

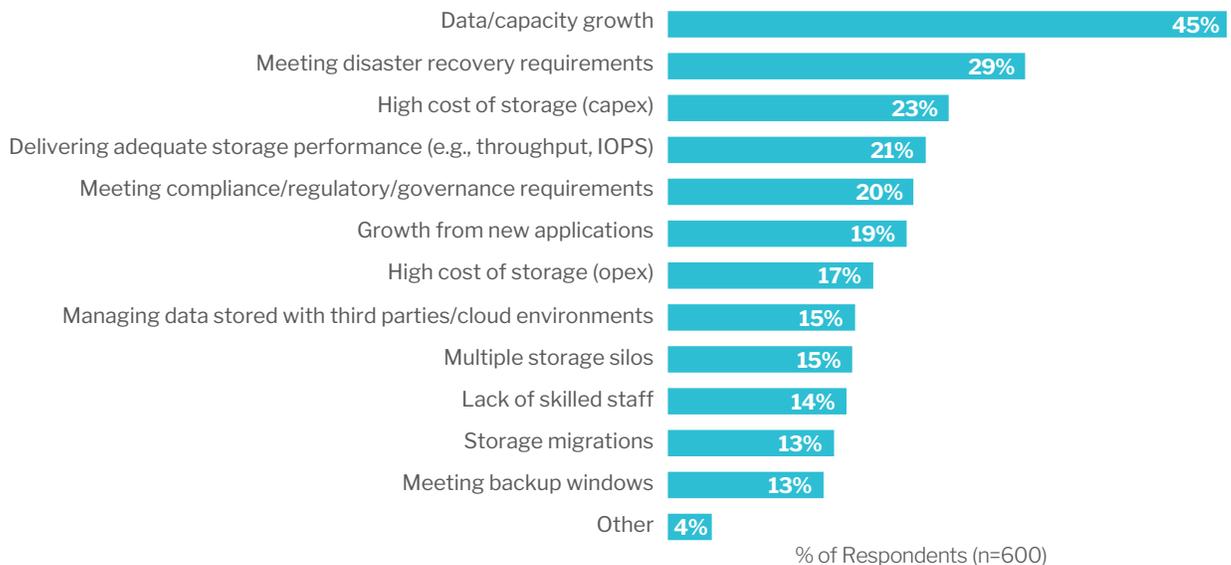
However, the recent emergence of the NVMe storage protocol challenges the disaggregation principle. The good news is that NVMe and its networked variant, NVMe Over Fabrics (NVMe-oF), allow even more workloads to enjoy the benefits of disaggregated storage. By boosting performance, they bring into the fold the minority of workloads that, for performance reasons, were previously restricted to non-disaggregated storage.

The bad news is that realizing the potential performance of NVMe in shared, disaggregated storage systems is not simple. Storage system controller latency overheads and the intense parallelism of NVMe are the major reasons for this. Overcoming these hurdles requires fresh thinking and innovation. A number of approaches have been taken by vendors selling storage systems that were purpose-designed to maximize NVMe performance, but the huge majority of them involve major compromises with respect to disaggregation. 451 Research believes that IT organizations should be very aware of this issue, because of its impact on agility, and both capital and operational costs, which are among the major storage-related pain points reported by enterprises (see figure).

Top Storage Pain Points

Source: 451 Research’s Voice of the Enterprise: Storage, Budgets & Outlook 2018

Q: What are your organization’s top pain points from a storage perspective? (Select up to 3.)



451 Research is a leading information technology research and advisory company focusing on technology innovation and market disruption. More than 100 analysts and consultants provide essential insight to more than 1,000 client organizations globally through a combination of syndicated research and data, advisory and go-to-market services, and live events. Founded in 2000 and headquartered in New York, 451 Research is a division of The 451 Group.

Business Impact

MODIFYING EXISTING STORAGE SYSTEMS TO USE NVME DOES NOT MAXIMIZE PERFORMANCE.

The majority of NVMe-powered disaggregated or SAN storage systems being sold today are still based on internal architectures that date from the disk era. Although the performance of these systems has been increased by moving from SAS to NVMe flash drives, the IO latencies imposed by their main x86 controllers prevent the full NVMe performance from being realized.

OFFLOADING PROCESSING ONTO HOST SERVERS DRIVES UP COSTS BY PREVENTING DISAGGREGATION.

One way to solve the aforementioned problem is to offload the storage controllers by moving some of the tasks they complete onto the application servers. This has been done in some purpose-designed NVMe storage systems, but it results in storage and compute no longer being disaggregated. It is very arguably a step backward in the evolution of datacenter storage, since it requires storage software agents to be maintained on application servers. That increases operational costs and complexity, and consumes host CPU cycles intended for application processing.

RELYING ON APPLICATIONS FOR DATA SERVICES ALSO DRIVES UP COMPLEXITY AND MANAGEMENT OVERHEADS.

Another approach is for purpose-designed NVMe storage systems to forgo data services such as snapshots or RAID protection against drive failures, and instead rely on applications running on host servers to provide those services. That can result in severely limited data services, and in all cases it forces IT staff to manage those services from within application silos – again preventing disaggregation, and driving up complexity and administration overheads and costs.

ARCHITECTING STORAGE TO MAXIMIZE NVMe PERFORMANCE WHILE RETAINING SERVICES MAINTAINS DISAGGREGATION.

Yet another approach is to create storage system architectures that avoid IO bottlenecks, while still providing storage and data services. Unlike many conventional storage systems that simply resemble dual-socket servers, such systems are based on innovative internal architectures that separate data IO from storage services. This allows the performance potential of NVMe to be achieved, while maintaining true disaggregation of storage and compute resources.

Looking Ahead

So far the disruption caused by the transition of datacenter storage from disk to flash has been limited, and the majority of all-flash storage systems currently being sold by incumbent suppliers are modified versions of systems that were designed in the disk era. That has been good for customers, and not just because it allows them to continue using products that have become integral parts of their infrastructure (which is why incumbent suppliers have taken this approach), but also because it has preserved disaggregation.

But the inevitable switch from the disk-era SAS and SATA storage protocols to the solid-state NVMe protocol will drive the industry to adopt new internal storage architectures. The industry has never stopped demanding higher processing performance, and currently the applications that are in most need of new levels of performance include, but are not restricted to, the hottest areas of IT – namely, machine learning and big-data analytics. The emergence of faster alternatives to flash as persistent storage, currently being spearheaded by Intel's Optane memory, will add to the pressure to develop new architectures because it will make even greater storage performance possible. The storage architectures that truly preserve disaggregation will enjoy major advantages.



With the latest version of its award-winning platform, Pavilion Data continues to enhance a design for NVMe that avoids IO bottlenecks, while providing storage management services like consistent snapshots and thin provisioning. By reducing latency from the host, across NVMe-oF networks and through the array to 40 microseconds, boosting read and write throughput of 90GB/sec and 120GB/sec, respectively, and adding SWARM Recovery for rapid node rebuilds, Pavilion Data is defining the playing field for NVMe-oF disaggregation.

To learn more, please visit www.paviliondata.com.