# WHAT QOS MEANS TO COMPUTING APPLICATIONS

Showcasing 5 Nines

## INTRODUCTION

In modern computing architectures it has been shown that being the 'fastest' is not always the best solution for maximizing performance of compute and storage infrastructure.

While we continue to see the shear IOPS performance of storage increasing with the move from SATA/SAS to NVMe as well as PCIe Gen 3.0 to Gen 4.0, it has been found that maximizing the bandwidth of these interfaces does not necessarily lead to an overall better performing system solution.

This is especially true on platforms now available that can support 48 drives or more in a 2U chassis, where we find that individual drive performance specs become irrelevant when the system level interfaces are incapable of aggregating that much raw bandwidth (BW). In effect, the drives go faster, but the system interface BW is still limited – you simply can't force any more data into the pipe.

Limited by CPU performance and/or networking interfaces, having that much sheer speed at the device level creates an asymptotic performance cap, and this bottleneck comes at a high cost. A storage device delivering a higher interface bandwidth requires an expensive storage subsystem using a high-power controller with many channels, additional memory and a large die size.

A new trend has been emerging that addresses this very problem. Whether it's called performance optimization or a

focus on Quality of Service (QoS), the approach is to have storage, like NVMe drives, that provide a high internal Read bandwidth with an optimized Write performance, tuned to maximize customer application use of many drives per system. With this focus on latency, each storage device has the ability to respond to host requests in an extremely predictable manner.

A common measure for this consistency is measured by a number of "nines". For a storage device to deliver a 5 nines value, it must complete tasks within a specific target window 99.999% of the time. Another way to describe this is that less than 1 in 10,000 commands completes outside a window of time regardless of the activities within that device.

### A TRAVEL-BASED EXAMPLE

An analogous way to look at this is via the model of boarding an airplane. The most efficient way to ensure success would be to board from the back, window seat first, then middle then aisle, progressively moving forward to the front. This would provide a boarding process that could achieve a 99.999% or 5 nines efficiency. But because the boarding process is by class and member status, it creates a totally random process and the efficiency drops below 1 nine of QoS.

### SOLVING COMPUTING APPLICATION QOS NEEDS

To provide this improved level of customer experience, a new approach to data storage is required, a solution that maximizes the Read experience and optimizes the Write performance, power consumption and scalability of the platform.

NGD Systems has developed a solution with patented Elastic FTL, Programmable ECC and data placement algorithms integrated into the controller architecture. The core of the platform manages Read vs Write ensuring house cleaning efforts do not negatively impact host QoS. When plotted, a narrower distribution of response times indicates higher consistency.

Figure 1 below shows the 5 nines distribution for the NGD Systems Intelligent Storage platform and that of a competitive solution running identical workloads. For this test, the drives were run against an FIO (Linux test platform) script with a customer workload that mimics their use case of mixed Reads and Writes.
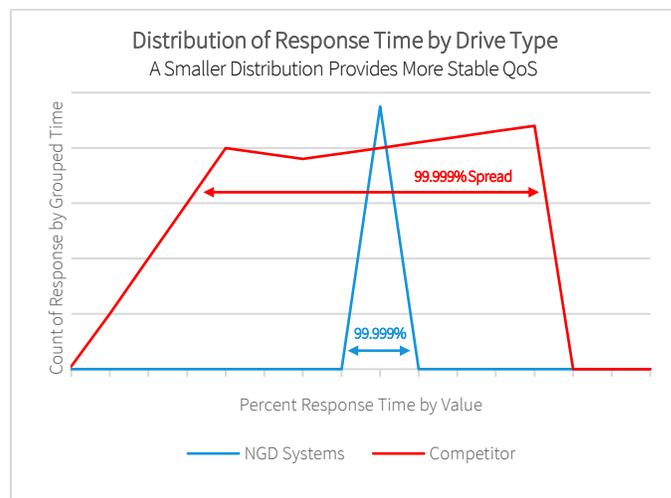


**FIGURE 1**
5 Nines Distribution for NGD Systems vs. Competitor
FIO Results Using Identical Mixed Reads and Writes Workload

The result: Optimized response time performance that scales across platforms and workloads in an NVMe drive that only consumes 12 watts of power.

### ABOUT NGD SYSTEMS

NGD Systems manufactures the world's most advanced computational storage drives (CSDs) and is fundamentally changing the IT Industry by bringing compute to data and achieving new levels of performance required by the next generation of data intensive applications. This enables customers to deliver results, like Hadoop acceleration or AI execution and address the most demanding service level objectives without requiring added hardware (GPU, FPGA, or other accelerating solutions). For more information, please visit https://www.ngdsystems.com.