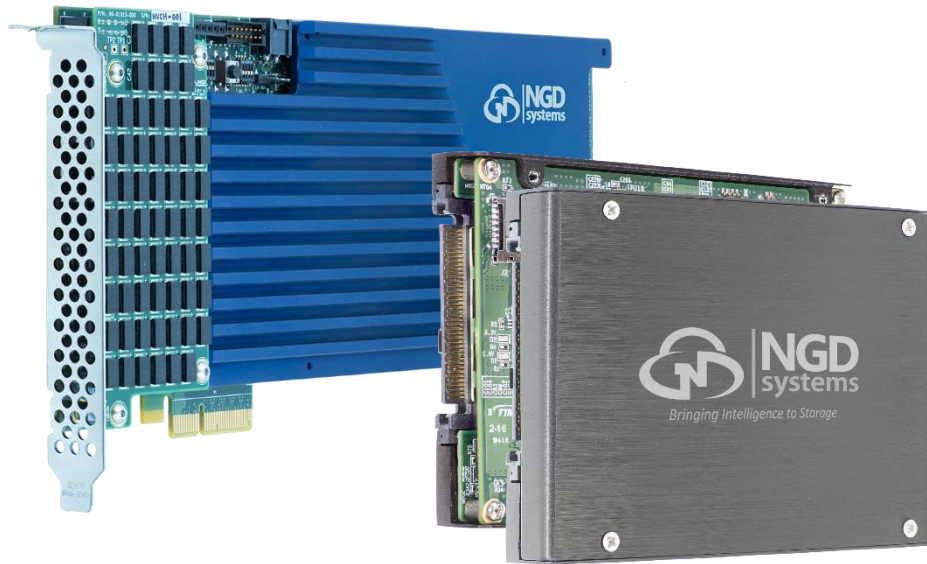




*Bringing Intelligence to Storage*



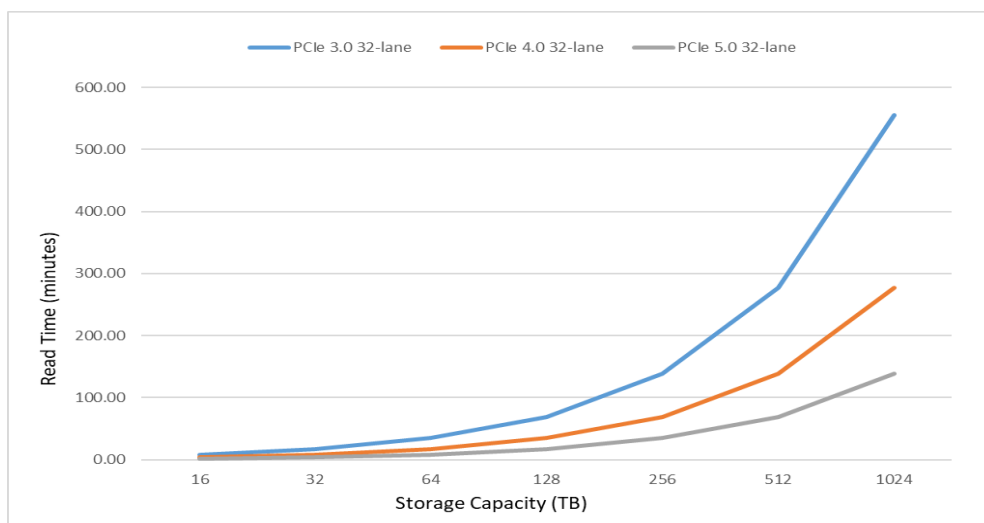
## **NGD Systems: Introduction to Computational Storage**

Updated: June 2018

**Executive Summary:** The advent of high-performance, high-capacity flash storage has changed the dynamics of the storage-compute relationship. Today, a handful of NVMe flash devices can easily saturate the PCIe bus complex of most servers. To address this mismatch, a new paradigm is required moving computing capabilities closer to the data. This concept, which is known as Computational Storage, provides storage platforms with significant compute capabilities within the storage device, reducing the computing demands on servers. For applications with large data stores and significant search, indexing, or pattern matching workloads, Computational Storage offers much quicker results than the today's traditional scenario of requiring data movement into memory and having the CPU scan the complete large data sets. This "In-Situ Processing" enables existing applications to scale compute within the storage devices making it possible manage data sets in place and reduce memory requirements and CPU load. Because computation capabilities scale linearly as storage is added into compute nodes, In-Situ Processing can enable new classes of applications for enterprises and cloud service providers.

## NGD Systems: Closing the Storage-Compute Gap

**Background:** The graph in Figure 1 below shows the read time **in minutes** necessary to read a given amount of data across various speed 32-lane PCI Express® (PCIe™) busses, from PCIe 3.0 to PCIe 5.0 (planned after 2020). As you can see on the chart, it takes over nine hours to read 1PB (1024TB) across a PCIe 3.0 bus and with a PCIe 5.0 bus, this time is still well over 2 hours. 1PB is not even a ‘normal’ amount of storage, noting that by 2019 EB-scale will be standard. Given this information, regardless of how ‘fast’ the storage interface becomes, the need to move data is becoming the largest bottleneck in all aspects of data management. There needs to be a move to a better way of getting analytics from data stored traditionally today.

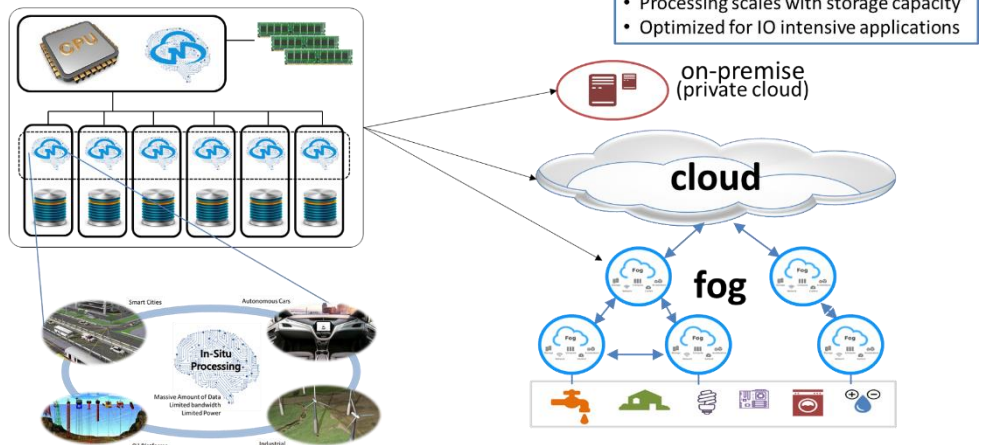


**Figure 1: PCIe Read Times by Storage Capacity**

The impact of this gap on problems with large dataset is immense. Many big data problems such as analyzing pattern/facial recognition image sets and searches of compressed datasets become out of the realm of economical real-time processing. The current solutions available in market are to take this data and split across many clusters. This conventional computing architectures significantly increases costs (both CapEx and OpEx), power consumption, and most important, the “physical footprint” of the system. In many cases, this also makes such problems out of the realm of reality or impossible to solve regardless of many added compute cores are present, since there is simply not enough ‘time’ to move all the data out of standard storage devices.

**The Solution – In-Situ Processing:** Stop Moving the data! Moving Compute processing into the storage devices accomplishes more capabilities and success compared to simply adding more compute away from the data. Moreover, if you could read the data at the speed of the SSD’s internal busses, you could reduce load times for each processor by 5X-10X when compared to the PCIe bus. Such an approach would certainly eliminate the “data wait” problem above. This is the concept of In-Situ Processing, is at the heart of NGD Systems patented Computational Storage solutions. In-situ processing is an extension of the “data locality” concept utilized distributed processing systems such as in Hadoop-based distributed big data solutions, where processing work is moved to where the relevant data is, rather than moving the data.

### In-Situ Processing

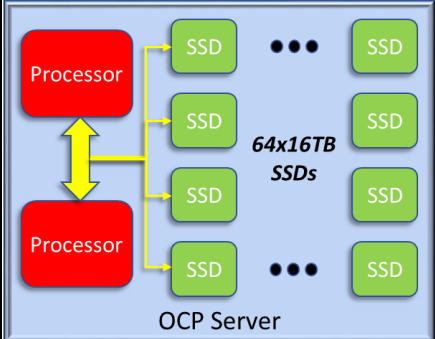
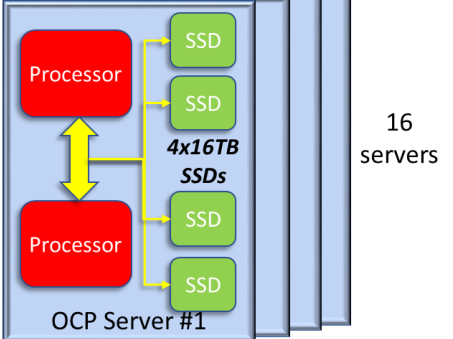
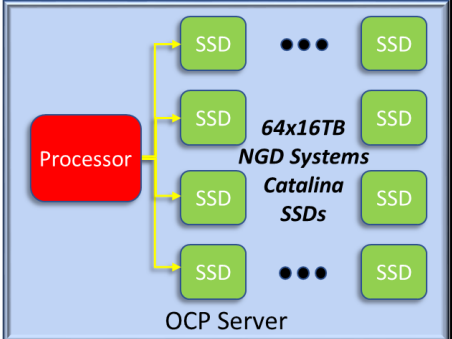


To understand how this concept changes the storage-compute dynamic, we will compare three configurations based on Open Compute Platform (OCP) servers, each with a 1PB dataset:

**Configuration 1:** One dual-processor OCP server with sixty-four standard 16TB U.2 SSDs.

**Configuration 2:** Sixteen (16) dual-processor OCP servers, each with four 16TB U.2 SSDs.

**Configuration 3:** One single-processor OCP server with sixty-four (64) U.2 NGD Systems Catalina 16TB drives, each consuming 12W.

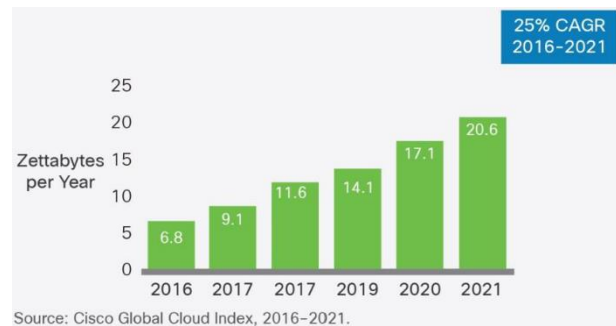
Configuration 1	Configuration 2	Configuration 3
 <p>OCP Server</p>	 <p>OCP Server #1</p> <p>16 servers</p>	 <p>OCP Server</p>
<p><b>One (1) dual-processor OCP server with 64 x 16TB standard SSDs</b></p> <p>Time to read 1PB dataset: 533 minutes</p> <p>Total power consumed: 1,370 W</p> <p><b>Total energy consumed: 12.17 kWh</b></p>	<p><b>Sixteen (16) dual-processor OCP servers with 4 x 16TB standard SSDs</b></p> <p>Time to read 1PB dataset: 36 minutes</p> <p>Total power consumed: 12,320 W</p> <p><b>Total energy consumed: 7.39 kWh</b></p>	<p><b>One (1) single-processor OCP server with 64 x 16TB NGD Systems Catalina SSDs</b></p> <p>Time to read 1PB dataset: 42 minutes</p> <p>Total power consumed: 1,348 W</p> <p><b>Total energy consumed: 0.94 kWh</b></p>

These configurations as shown on the previous page, showcase the time to access an entire dataset (1PB), along with the total power and total energy consumed by each configuration. Clearly, the access time and energy consumption improvements from the one dual-processor OCP server with standard SSDs to the one single-processor server and NGD Systems Computational Storage devices is considerable - it is an order of magnitude better on the metrics. Meanwhile the 16-server OCP cluster configuration provides a slightly better access time at the cost of higher power consumption. The 16-server cluster consumes a massive 12KW of power, and over 7 kWh of energy, and occupies a full rack for such a small data set. This configuration, while possible is highly impractical from a power, cooling, CapEx and OpEx standpoint.

**Use Cases:** There are a variety of use cases which lend themselves to in-situ processing. In general, these use cases are read-intensive with large numbers of parallel operations that are performed on the dataset, especially pattern matching, indexing, and searching. We will explore a few here where In-Situ Processing enables new levels of performance and/or capabilities.

**Hyperscale Storage:** By 2020, roughly 485 hyperscale data centers (HDC) will contain 57% of all data stored in datacenters worldwide. Moreover, the traffic measured in ZB (1000 EB), shown below within these HDC will quintuple by 2021<sup>2</sup>. For a hyperscale storage architect, finding a way to avoid moving data from device to device even in the same rack will become a paramount need, as will reducing the power consumed to process this data.

In-Situ Processing offers the best approach today to mitigate both the data movement and power consumption issues. More importantly, In-Situ Processing significantly improves execution times for applications that make use of operations such as indexing, parallel searches, and pattern matching – all of which are becoming more and more common as AI/ML applications proliferate.



**Machine Learning/Embedded Artificial Intelligence:** Artificial Intelligence (AI) has achieved very promising results in areas such as computer vision (CV), speech recognition, and natural language processing. For example, in a smart city data such as videos or images captured from many distributed cameras need to be automatically processed using video analytics; i.e., object detection, object tracking, facial recognition, image classification, and scene labeling. Vast amounts of cloud data can be used for training on powerful platforms to create generalized yet accurate models.



However, there is an opportunity to move the inference, and in some cases the learning, onto “Intelligent” storage. This is where one can take advantage of a distributed system composed of the computing resources available in a multitude of Computational Storage devices to implement shallow learning with weightless neural networks.

**Intelligent Edge Computing for IoT (OpenFog):** The Internet of Things (IoT) aspires to collect mountains of data by instrumenting nearly everything in our existence. This creates challenges both in storing and processing this data. Gartner Research expects that there will be 20.4 billion IoT devices connected to the internet by 2020, generating 5X the data that we generate today. Cisco expects that IoT devices will generate 403 zettabytes/year by 2018<sup>2</sup>. A single connected aircraft can generate 40TB of data per day<sup>2</sup>, while an autonomous car may generate 2TB of data per hour, according to Intel<sup>3</sup>. Clearly, this data cannot be simply “sent home” over the network to datacenters, even if they are close in the cloud. The key to making constructive use of this IoT data deluge is the ability to process it “In-Situ” or in place, which is exactly what Computational Storage provides for IoT. By In-Situ Processing this data, IoT nodes can make intelligent updates to the cloud without the burden of the power, space, or cost of a server.



**Leading-Edge Intelligent Storage Solutions:** At NGD Systems, we are creating the paradigm shift products of Computational Storage with In-Situ Processing for the way to make the future successful. Our patented architecture and available Catalina-2 products make the deployment of many large dataset applications possible and practical, whether from an access time, power/cooling, or real estate.

If you would like to find out more on how you can benefit, please contact us for further discussions on how Computational Storage can help solve your data center and business issues.

[Info@NGDSystems.com](mailto:Info@NGDSystems.com)

[Sales@NGDSystems.com](mailto:Sales@NGDSystems.com)

+ 1-949-870-9148

**Footnotes:**

1 – [Gartner Press Release, Feb 7, 2017](#)

2 – [Cisco Global Cloud Index: Forecast and Methodology, 2015-2020 \(2016\)](#).

3 – [Patrick Nelson, Network World, Dec 7, 2016](#)