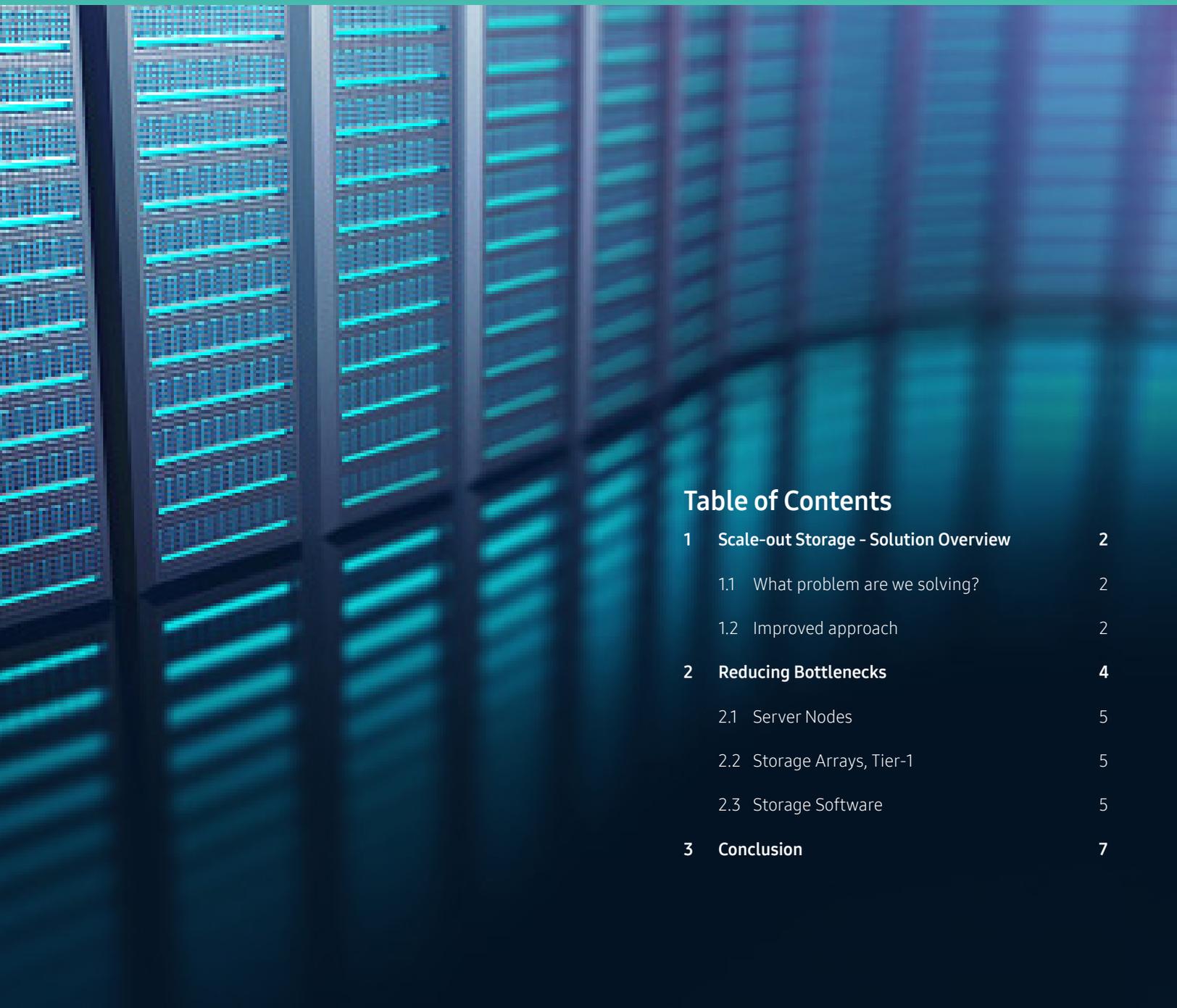# Optimizing performance in HPC/AI systems

A new reference design from Samsung and Applied Data Systems

Principal Author: Hubbert Smith | Samsung Semiconductor, Inc.

## Table of Contents

# 1  Scale-out Storage – Solution Overview

## 1.1  What problem are we solving?

Modern-day High Performance Computing (HPC) requires system designs that balance compute, networking and storage to eliminate bottlenecks which can hinder time-to-results for high-value, time-sensitive tasks. For well over a decade, HPC decision-makers have faced rising storage costs and rapidly growing demand for storage capacity. Storage growth, without the benefit of improved performance, has created a significant bottleneck. And so, HPC data centers are seeking disruptive new solutions that:

- Offer a balanced combination of resources – Primary Storage, Secondary Storage, Network, DRAM, CPU and GPU.

- Meet the tremendous need for higher productivity – reduced time to outcome, more jobs per day/increased volume of data per job.

- Enable less repetitive user interaction – reduced 'brute force' data movements such as 'cp' and 'rsync'.
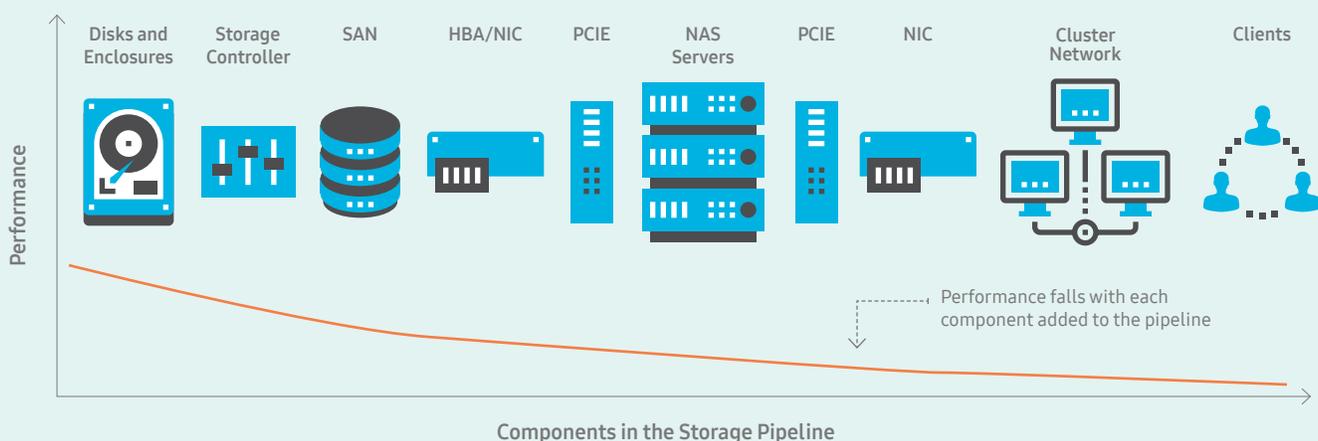
## 1.2  Improved approach

Previous approaches have stagnated in the face of multiple bottleneck issues:

- Network roadblock, often at 1Gb/sec or 10Gb/sec

- Low CPU core count, primarily with Xeon E5-2600 family, 4-core

- Reaching the limits of DDR3 and DDR4 memory, typically 128GB or 256GB

- Log-jams created by storage based on hard disk drives (HDD)

OEMs have encountered significant hurdles in maximizing storage efficiency when using servers with directly attached HDDs, or scale-out storage software, which usually have not been optimized for performance. Typically, previous solutions employed the "then-current" set-up with CPU, HDD and software.

### Figure 1. The Old Pipeline



Performance falls with each component added to the pipeline

Performance
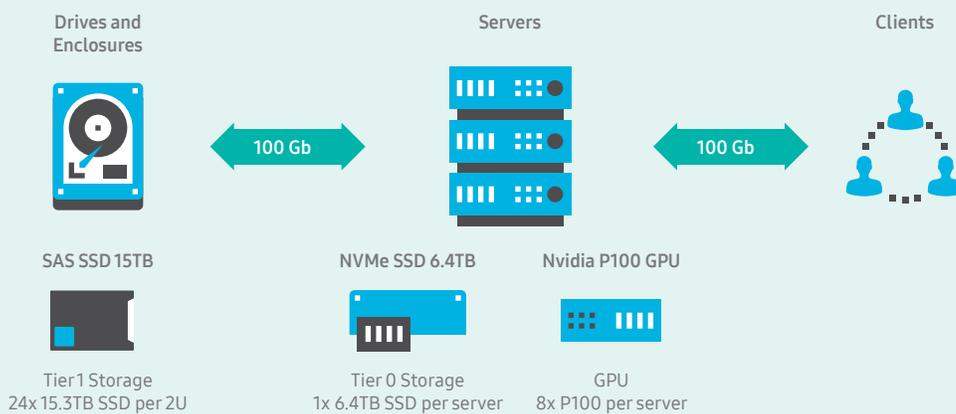
Components in the Storage Pipeline

# 1  Scale-out Storage – Solution Overview (continued)

However, the world has changed. Today, networking is significantly faster, CPUs have more cores, DRAM is more efficient, SSDs have replaced HDDs in most instances (except for cool storage), and storage software has improved considerably.

In an increasing number of HPC storage scenarios, architectures are using GPU-based systems that require higher IOPs as well as greater bandwidth.

## Figure 2. The New Pipeline



Drives and Enclosures

Servers

Clients

100 Gb

100 Gb

SAS SSD 15TB

NVMe SSD 6.4TB

Nvidia P100 GPU

Tier 1 Storage
24x 15.3TB SSD per 2U

Tier 0 Storage
1x 6.4TB SSD per server

GPU
8x P100 per server

**Applied**
DataSystems

**SAMSUNG**

# 2  Reducing Bottlenecks

Let's look at a new generation reference design for HPC clusters, from a system-wide purview. Applied Data Systems and Samsung have collaborated to provide what is arguably the most efficient way to process data in an HPC environment.

First, this solution is designed to avoid performance bottlenecks prevalent in previous generations of HPC. Second, it has been specifically designed accommodate GPUs used to expand capabilities for artificial intelligence and machine learning. Lastly, this HPC storage reference design includes a balanced deep-storage process for collecting/capturing high volumes of information, feeding the data into a growing arsenal of AI resources, then capturing the resulting AI-driven data and protecting it.

Time-to-solution is important, and therefore, performance is essential for high-value, time-sensitive tasks such as energy exploration, genomic sequencing, stock trading and fraud analysis.

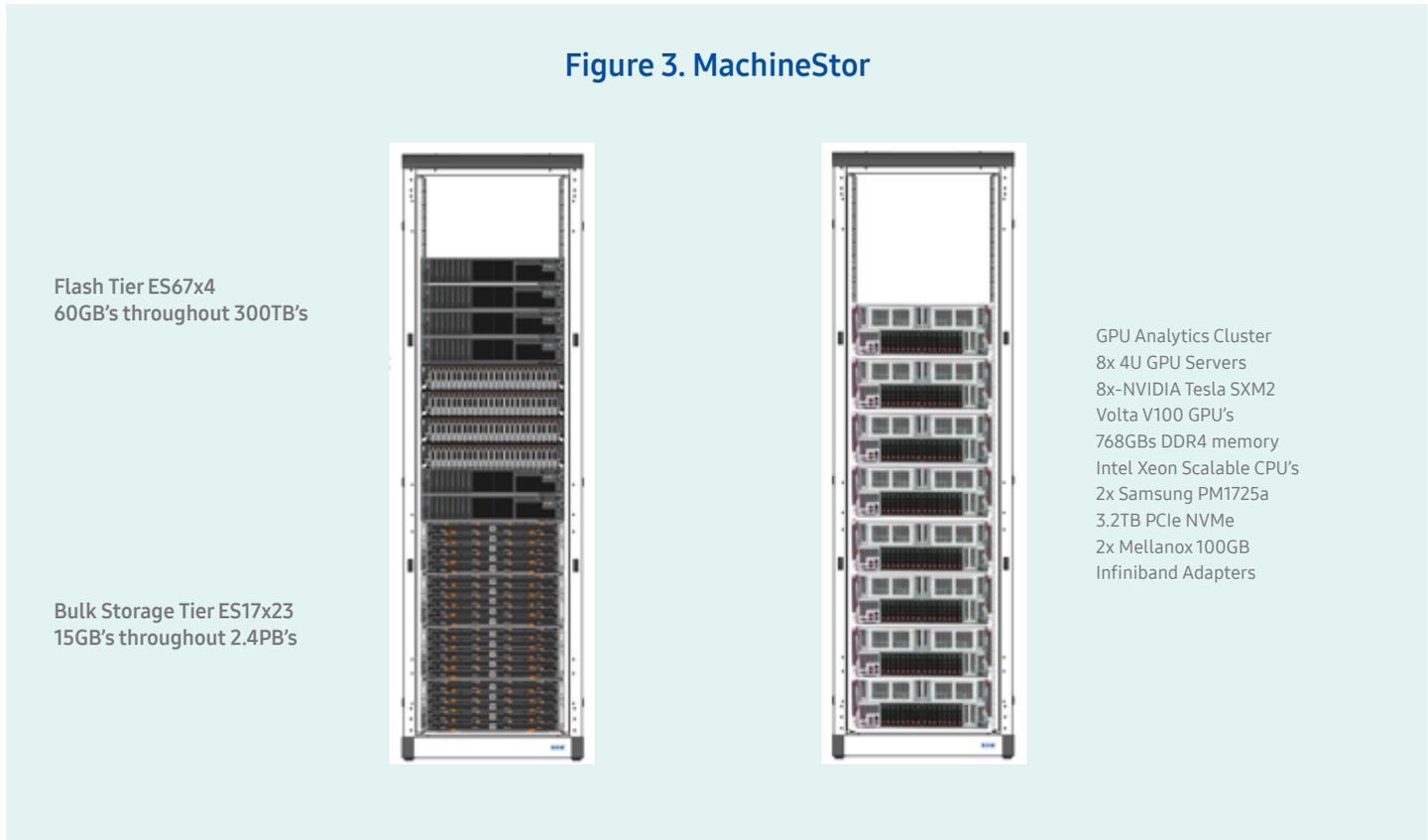The sooner HPC results are delivered, the sooner the desired financial, or life-saving, outcomes are realized.

This reference design is focused on delivering results quicker, fully accepting that:

- Networks now are typically 100Gb Ethernet or 100Gb Infiniband

- There is greater use of GPU-based computing nodes

- The CPU core count is rising, currently up to 28 cores

- Systems can accommodate up to 3TB of DDR4 – with 128GB per DIMM, and up to 24 DDR4 DIMM slots

- Storage uses NVMe for Tier0 and SAS SSD for Tier1

- Storage software has improved dramatically, using NVMe over Fabrics for scale-out block with a CPU bypass. Users now scale filesystems in parallel such as IBM Spectrum Scale or Lustre, creating a scalable global namespace that is easy to manage (no more 'cp" or 'rsync')

| Server CPU, DRAM | Servers with 2 Intel E5 v4 series processors<br>512GB DDR4 – Samsung |
|---|---|
| Server GPU | Nvidia P100 GPU, V100 GPU<br>x8 per server |
| Server Network HBA | Mellanox ConnectX-5, 100Gb/s NICs with dual-ports<br>MCX 515A-CCAT<br>Mellanox OFED 4.0-2.0.0.1 driver<br>Link here for driver G-zip, Must support RoCE V2 |
| Tier-0 SSD in Server | Half-Height, Half-Length (HHHL) PCIe card<br>Samsung 1725a 6.4TB SSD |
| OS Version | Linux 4.4.0-040400-generic |

| Tier-1 Storage Array | NetApp EF-Series 560, 570 or equivalent |
|---|---|
| Storage Devices | 24x 15.3 SAS SSD, from NetApp |

| Software Scale-out Block | Excelero NVMesh 1.2 Link here |
|---|---|
| Software<br>Scale-out File | IBM Spectrum Scale<br>(formerly IBM General Parallel File System – GPFS) |

# 2 Reducing Bottlenecks (continued)

## Figure 3. MachineStor

**Flash Tier ES67x4**
**60GB's throughout 300TB's**

**GPU Analytics Cluster**
**8x 4U GPU Servers**
**8x-NVIDIA Tesla SXM2**
**Volta V100 GPU's**
**768GBs DDR4 memory**
**Intel Xeon Scalable CPU's**
**2x Samsung PM1725a**
**3.2TB PCIe NVMe**
**2x Mellanox 100GB**
**Infiniband Adapters**

**Bulk Storage Tier ES17x23**
**15GB's throughout 2.4PB's**

## 2.1 Server Nodes

These nodes include 8x Nvidia P100/V100 GPU. They typically consume approximately 500MB/s – 700MB/s of bandwidth per GPU.

Eight GPUs will, in aggregate, generate 4GB/s – 5.6GB/s of demand on Tier-0 storage. The design specifies Tier-0 as a 6.4TB PM1725a HHHL – which delivers up to 6GB/s. The Samsung PM1725a provides a great ratio match with an 8x P100/V100 GPU.

## 2.2 Storage Arrays, Tier-1

The Tier-1 storage arrays have been optimized for high performance, low latency and high capacity. In fact, a 24x 15.3TB SSD yields 367TB of raw capacity.
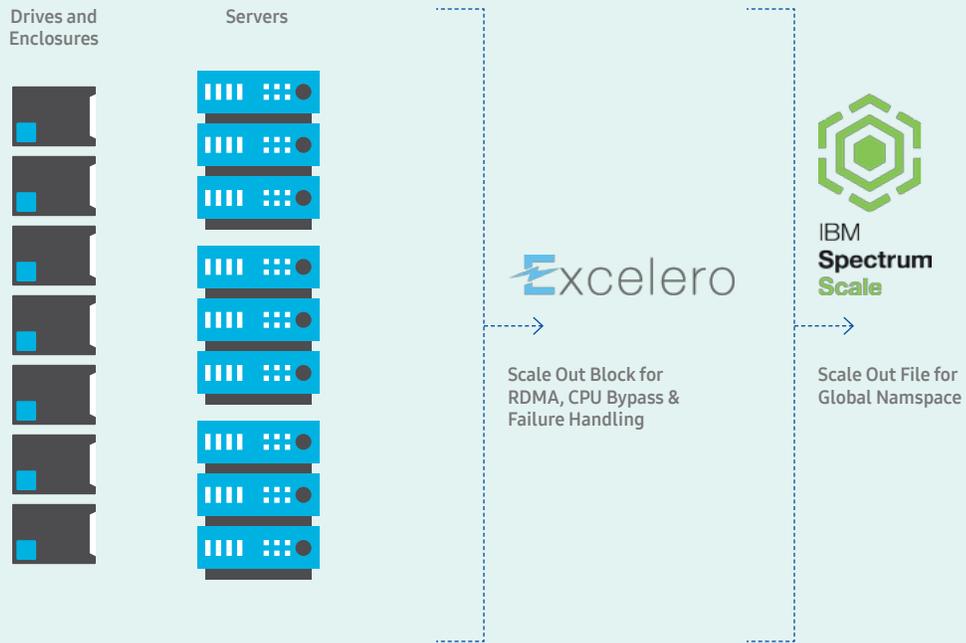
## 2.3 Storage Software

The Applied Data - Samsung design applies Excelero NVMesh. NVMesh provides a global block space, with CPU bypass, for significantly increased performance.

This highly innovative design employs IBM Spectrum Scale (formerly IBM General Parallel File System) which provides a global file namespace. The design results in a significantly improved operational model, and eliminates the need for brute force data commands like 'cp' or 'rsync'. It greatly improves operational efficiencies and reduces the potential for human error.

Applied
DataSystems

SAMSUNG

# 2  Reducing Bottlenecks (continued)

### Figure 4. Software Defined Storage

Drives and Enclosures

Servers

Excelero

IBM
**Spectrum**
**Scale**

Scale Out Block for RDMA, CPU Bypass & Failure Handling

Scale Out File for Global Namspace

# 3  Conclusion

From extensive industry experience, we know that important scientific, medical, engineering, finance and exploration problems are ever-evolving. We know these problems have evolved past what linear programming can overcome, so AI and ML are now essential. We also know that the volumes of data being stored and orchestrated are exponentially growing, so balanced and de-bottlenecked IO is essential for orchestrating all of the data flowing in and out of these HPC/AI clusters.

We encourage you to pursue your own evaluation soon, leading to a Proof of Concept.

Applied Data Systems, Inc. (ADS) is focused on delivering engineered High Performance Computing (HPC) solutions for customers.  Applied Data Systems specializes in optimizing high-speed parallel file systems, designing high-speed networks, building HPC clusters, and customizing servers and storage to provide as much performance and efficiency as possible.  With over 50 years of combined HPC-specific experience, Applied Data Systems' core team has designed solutions for customers ranging from the largest HPC labs in the world to the Fortune 500.

## For assistance, please contact:

**Applied Data Systems**

12180 Dearborn Place
Poway, CA 92064
(844) 371- 4949
info@applieddatasystems.com

OR

**Samsung Semiconductor, Inc.**

Hubbert Smith
hubbert.s@samsung.com

## For additional info:

https://www-03.ibm.com/systems/storage/spectrum/scale/big-data-storage/
http://www.samsung.com/semiconductor/products/flash-storage/enterprise-ssd/
ssd@samsung.com

Applied DataSystems

SAMSUNG