

# **Accelerating Real-Time Big Data**

Breaking Through Performance and Scale Out Barriers – A Storage Solution for Today's Hot Scale Out Applications



## Agenda

- Everything related to storage is changing!
  - The 3rd Platform
  - NVM Express architected for solid state storage
- The Migration of Storage Intelligence
  - Storage Aware, Scale-Out Applications
  - NVMe controllers
- Storage Array Controllers are DEAD!
- Captive storage limitations
- Apeiron architecture NVMe over Ethernet
  - Noe / NVMe over Fabrics comparison
- The Apeiron Storage Solution
  - Captive storage vs. external storage performance
  - 18.4M IOPs in 2 rack units
  - The universal NVMe storage platform
  - Advanced features

# The 3rd Platform: A Fundamental Shift in the IT Industry



## NVMe<sup>™</sup> Delivers Higher IOPs and Better QoS



#### NVMe<sup>™</sup> delivers 18 µs average and 40 µs 99.99% interface latency. Other interfaces have outliers in 100s of µs as interface reaches saturation.

Results measured by Intel based on the following configurations. Intel Server Board S2600WTT with 28 E5-2695 CPUs, 2 sockets, 2.3 GHz clock speed per CPU, Ubuntu\* 14.04.1 LTS (GNU/Linux\* 3.16.0-rc7tickles x86\_64), idle=poll kernel settings, SAS HBA is LSI SAS9207-4i4e with controller LSI SAS 2308. SATA SSDs are Intel® SSD DC 3500 at 800 GB. NVMe SSD is Intel SSD P3700 at 1.6 TB. Workload details are Workload: 4K Random Reads using FIO – 4 + threads. Drives tested empty to test interface only (no NVM access.)

9



## **NVMe delivers**

- Performance
- Managability
- Robust ecosystem
- Well defined standard SSD form factor
- Steep innovation and healthy competition
  - Performance, durability, capacity and cost





# **The Migration of Storage Intelligence**

3rd Platform Storage

NVMe

2nd Platform Storage

1<sup>st</sup> Platform Storage

and control

•

•

Millions of developers (open source)

Storage aware applications (and OS)

In-memory data base, native tiering Sever is the critical component

Architected for Scale out - not scale up

Very High Performance persistent storage

Flash now, Storage Class Memory soon

> 500K lines of code in an NVMe controller

Simple storage drivers, SCSI, smarter devices

> 25M lines of code in storage controller SW release

**Direct Connect, Dumb Storage Hardware** 

Software (OS) centric storage management

6

Direct attached storage >= networked

Very intelligent storage devices

Network attached storage (SAN)

Array Controller centric intelligence

#### **THE 3rd PLATFORM**

Defining the integration and intersection of mobile, cloud, social, and big data



## Array Controllers are DEAD! (well, dying)

#### • Storage aware, scale-out, real time analytics

- Examples: Ad Tech, fraud detection, facial recognition, personalized user experience, enterprise big data analysis, etc.
- Application manages data placement across compute cluster
- Server is now the critical component HA / failover strategies required
- Application manages HA, data tiering and migration
- Complex, array centric storage SW = slow perf new apps just don't use it
- Multiple tiers of storage are a given and managed by the app.
  - DRAM, captive NVMe, captive SATA, external flash then maybe HDD
- NVMe storage devices
  - Solid state drives are much more reliable than HDDs
    - HDD MTBF drove the development of Storage Array Controllers (NetApp, EMC, HP . . .)
  - But flash wears out requires complex management code including loads of data movement "behind the curtains" in the device
  - NVMe standard was written for flash controllers with ample processing power
    - Excellent device management and monitoring
  - NVMe controller handles data movement
    - Turns the SCSI model upside down

Array Controllers add no value for Scale-out Apps

#### 2nd Platform Storage including all flash arrays



#### **3rd Platform External Storage** Direct Attach scale-out storage

**Applications** 

<u>Uşer Spac</u> Kernel ADS Mgnt Agent

**ADS Driver** 

- Software / application defined storage
- One, very high performance storage network
- Designed for scale out integrated switches
- Scales to 100s of servers, multiple petabytes



# Why Not Captive Storage



- Captive (direct attached) Storage
  - Limited total capacity and performance
  - No dynamic scaling
  - No SSD virtualization
  - No high performance data sharing / tiering across cluster
  - A severe management challenge
  - Inefficient power, cooling, rack space
  - Storage provisioning is tied to CPU scale out
  - PCIe board solutions are worse!

### Get the Storage Out of the Server!! (again) :

## Storage Network Protocol NVMe over Ethernet (NoE)

- Hardware accelerated, hardened Layer 2 Ethernet fabric
  Layer 3 robustness without the overhead and latency
- Fully integrated NVMe fabric (no external switching)
- The industry's lowest latency transport protocol delivers predictable performance at scale



NVMe Virtualization Software

L2 Ethernet Headers + 4B NoE Header/encapsulated PCIe TLPs

### **NVMeoF (RDMA) / NoE Comparison**



The standard is not tied to any particular physical layer RDMA approach adds between 26B and 96B of headers, in addition to NVMe Encapsulation

Flexible but adds complexity, link consumption and latency !

# **NoE / NVMe over Fabrics comparison**

NVMe over Ethernet (NoE)	NVMe over Fabrics
Transports NVMe commands	Transports Data (RDMA)
Optimized for Ethernet (minimizes overhead)	Transport independent (more complexity)
4 Byte per packet added overhead	>> overhead (depends on implementation)
Optimized for scale-out clusters	Architected for traditional storage arrays
Supports ANY standard NVMe SSD	????
Next gen (3D XPoint) ready	???? (demonstrated latency is a problem)
Shipping today	NVMeoF standard now approved

Apeiron is in production *today* shipping the highest performance scale out NVMe storage solution in the world

All mentioned brand names are registered trademarks and property of their respective owners. "3D XPoint is a trademark of Intel Corporation in the U.S. and/or other countries"



#### ADS1000 Scale-out NVMe Solution Unmatched Performance, Scalability and Efficiency

#### **NVM** EXPRESS

**C** Fabric Ports

24 NVMe 2.5" SSD



Serviceability

#### ADS1000 Performance (2U)

Capacity	38/76/154/192TB
Latency (NAND LIMITATION)	100us
Protocol Overhead	<3us (roundtrip)
Bandwidth sustained	72 GB/s
Random 4K reads	18.4 M IOPS



& Cooling Modules



#### **A New Standard in Storage Networking Performance**



#### 24 NVMe 2.5" SSD



The ultra low latency Apeiron network technology is 100% transparent to the servers\*

Apeiron vLUN's enable workload optimization across multiple NVMe drives

\* Please see the March 2016 ESG Whitepaper-"Validation of Apeiron Performance" at apeirondata.com

## **FIO Performance Benchmarks**



- The balanced architecture of the ADS1000 improves server utilization up to 3x (limited by SSD type)
- Apeiron's driver technology moves the bottleneck from the network to the NAND architecture itself (18.4M IOPs single ADS1000)
- Performance scales linearly to 100's of millions of IOPs as the system solution grows



Scaling serverBiv<sup>™</sup> & 6.4TB SSD's 150 TB → > 1.2 PB raw storage 15 M → 120 M IOPS, 4kB Random Read 64 GB/s → > 500 GB/s Bandwidth

Scaling serverA<sup>#</sup> 3.2TB SSD's 76 TB →>380 TB raw storage 15 M →>75 M IOPS, 4kB Random Read 64 GB/s → 320 GB/s Bandwidth

Scaling Internal SSD's (DAS)



## The World's Only Universal NVMe Platform

- Unlike captive storage, Apeiron enables independent scaling of servers and storage
- Compatible with ANY commercial NVMe drive-Data resides on appropriate SSD type for its value (Including 3D XPoint<sup>™</sup> technology)
- Adoption of NVMe SSD's is rapidly increasing; Only Apeiron can provide compatibility with all suppliers and drive profiles



The roadmap for density and performance of NVMe SSD's is accelerating; Apeiron passes this advantage to the customer

All mentioned brand names are registered trademarks and property of their respective owners. "3D XPoint is a trademark of Intel Corporation in the U.S. and/or other countries"



#### **NVMe Solution Comparison**



	System A	ADS1000
Rack Units	5U	4U (2x 2U)
Bandwidth	100GB/s	144GB/s
IOPs	10M IOPs	37M IOPs
SSD	Proprietary	Any SFF NVMe SSD
Latency	100us (avg)	100us
Interconnect	PCIe 3.0	40Gb Ethernet
Maximum Capacity	144 TB	4.6 PB, 9.2 PB Q3'16 (60 enclosures)
Intel 3D XPoint <sup>™</sup>	No	Yes 3D Xpoint = 7us latency*
Entry Level List Price	3x Apeiron	1/3 System A



#### Performance, Scalability and Simplification

"All the simplicity and promise of direct attached storage with the capabilities of network attached storage."

-Ahmed Shihab, VP Engineering, Amazon Web Services

