# Optimizing SSD for Database Applications

August 24th, 2017

Dr. Xueshi YANG
CEO, Shannon Systems

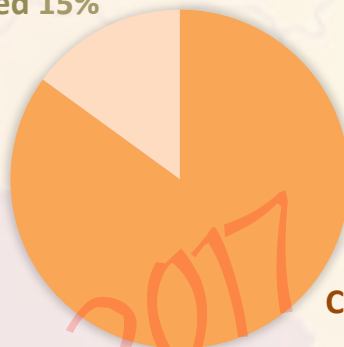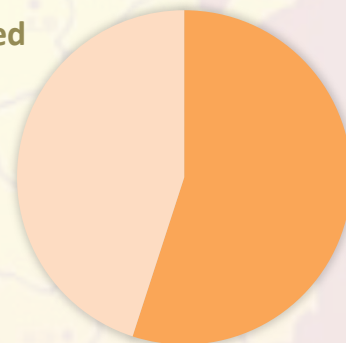Shannon Systems
宝存科技

Not covered 15%

Covered 85%

Shannon SSDs inside ~85% of internet companies

Not covered 45%

Covered 55%

We power more than 50% of e-commerce databases

or find your significant other...
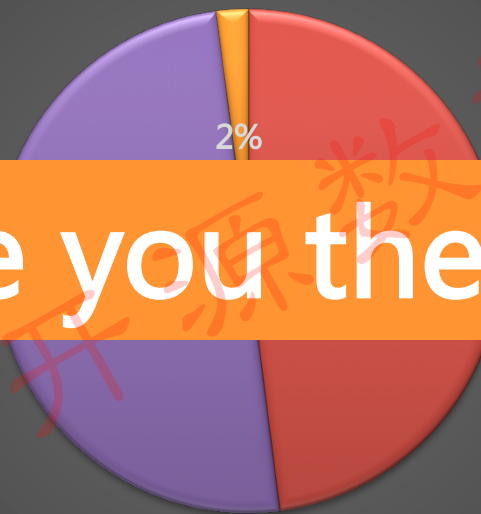
# 3D NAND Flash

- 96 layer NAND
- 4 bits/cell

**512Gb per die
1TB per single BGA package**

Source: WDC



1.5x layers

512G, 64L BiCS3

35% die shrink

512G, 96L BiCS4

# 5$^{Th}$ generation Shannon Direct-IO controller



## Tailor-made FFSA controller

- 1 Million IOPS
- 36TB capacity
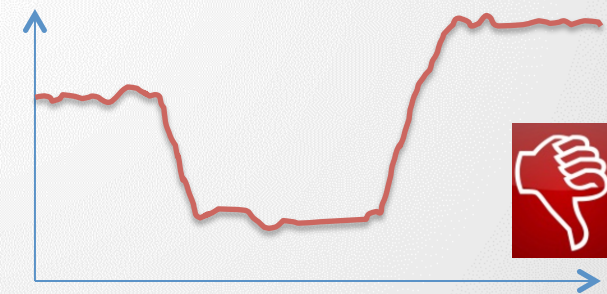- micro-second access latency

# Sharing by multiple Databases

- Capacity and performance sharing

# IO consistency

# IO separation

Approach: IOPS limitation per instance

# IO separation

- Limit the maximum available resources for each applications

# IO separation via streams



Source: Landsman,FMS2017

# IO separation via streams



Source: Landsman,FMS2017

# Networked SSDs

- IO access across nodes or space is the key enabling component



Multiple-Nodes Hyper-Converged Infrastructure for XenServer

# Networked IO

- RDMA technology enables SSD sharing across nodes with minimum latency overhead

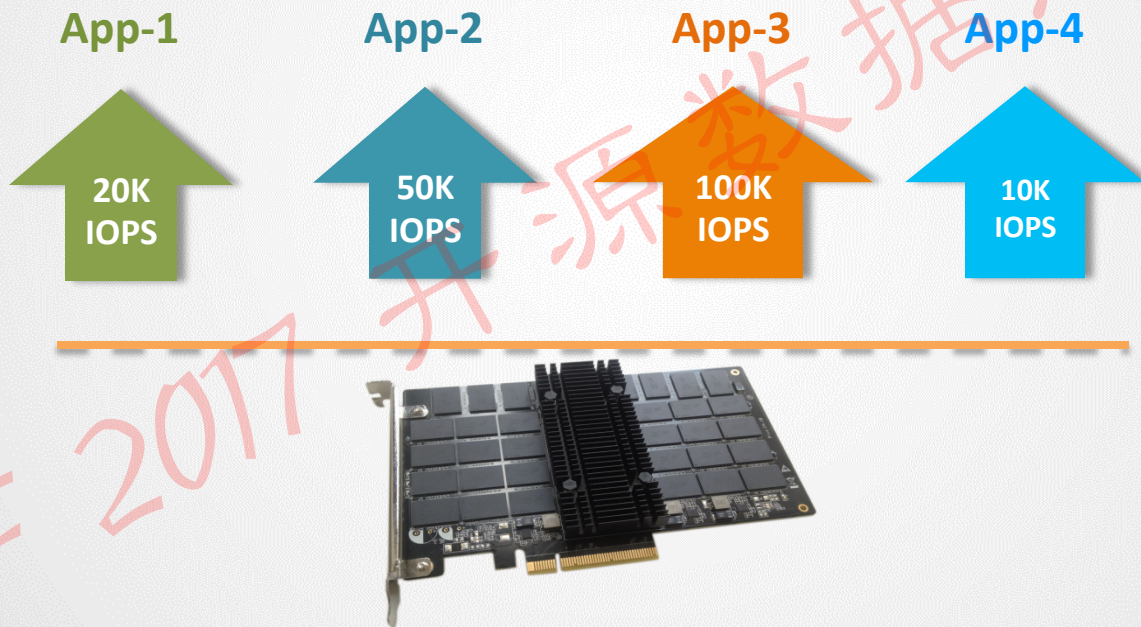    - e.g. NVMeoF

- Proprietary software layer ensures IO determinism across space

    - Access latency becomes node agnostic

    - Storage pool is constructed across multiple nodes

    - SSD pool becomes HA (highly available)

# Networked IO

- IO access through high speed network via RDMA



Access latency (us)



Random IOPS (K IOPS)

# IO consistency –RW separation

- Maintenance window to guarantee read access latency

| Deterministic read window | Maintenanc e window | Deterministic read window | Maintenanc e window |
|---|---|---|---|

GC/WL, background scrub etc

Shannon Systems
宝 存 科 技

# IO consistency –RW separation

- Pool multiple SSDs together via RDMA network

| Maintenance window | Deterministic read window | Maintenance window | Deterministic read window | Maintenance window |
|---|---|---|---|---|

| Deterministic read window | Maintenance window | Deterministic read window | Maintenance window | Deterministic read window |
|---|---|---|---|---|

| Deterministic read window | Maintenance window | Deterministic read window | Maintenance window | Deterministic read window |
|---|---|---|---|---|

**At any time, one copy exists for deterministic read IO**

Shannon Systems
宝存科技

Shannon Systems
宝存科技

# IO prioritization

- Provide priority service for certain traffic

# IO prioritization

- Higher priority processing for IO from a particular application
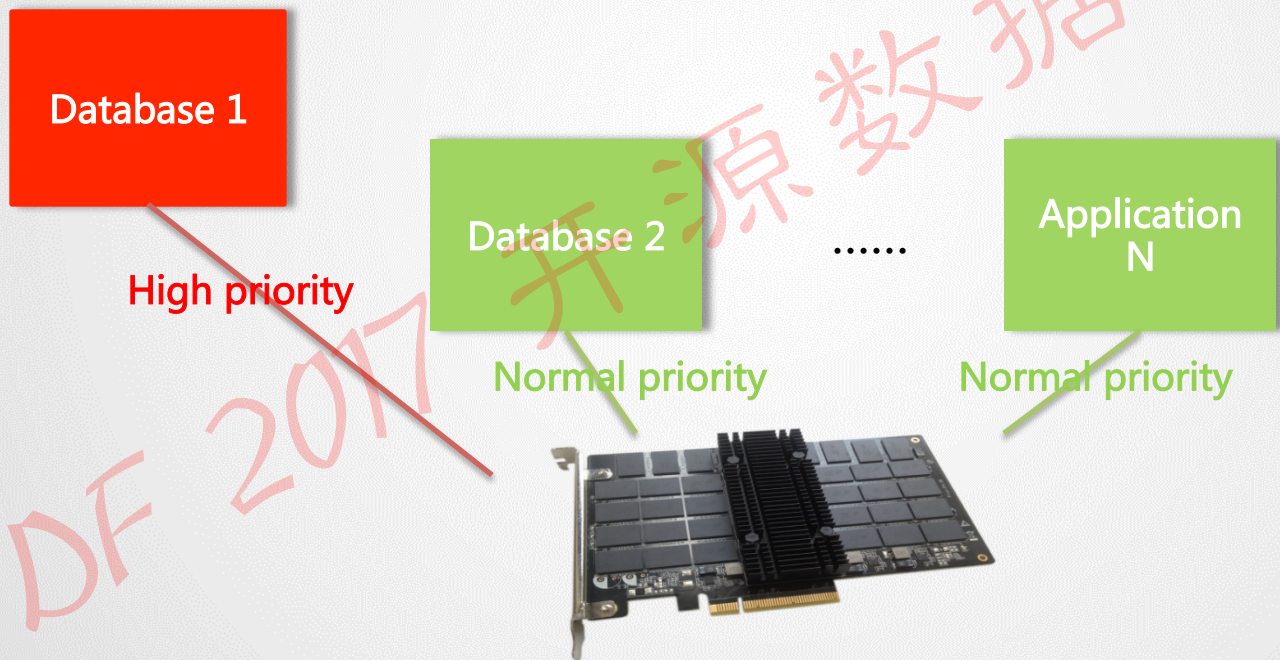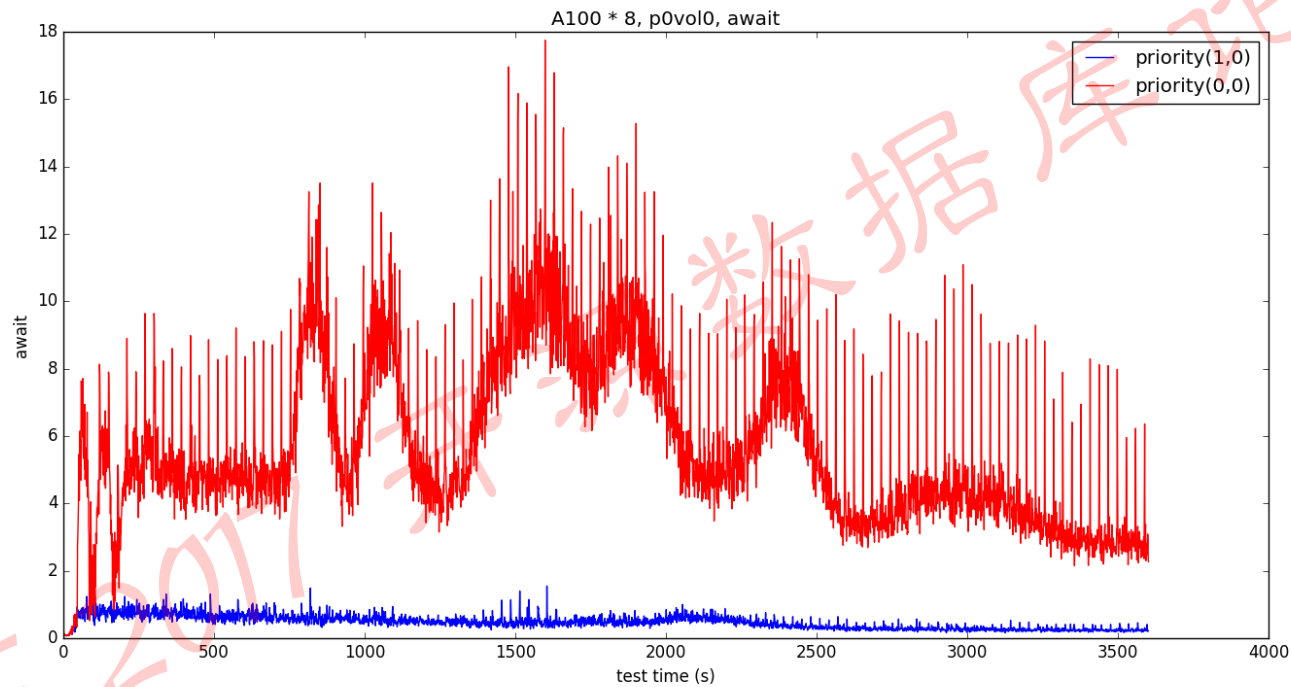- Priority Read/write operations

# IO atomicity

- IO (minimum size) sent to storage system
    - Success
    - Complete fail
- No intermediate state is allowed
- For example, Shannon Direct-IO PCIe Flash guarantees atomicity for any IO less than 32KB
- MariaDB supports atomic writes without using double write for improved performance.

# Summary

**1**

SSD capacity/ performance admits resource sharing among applications

**2**

IO determinism elevates QoS
- separation
- prioritization
- atomicity

**3**

IO determinism is enabled via software-defined SSD
- Intelligence is offered in host SSD driver

**4**

Database applications can benefit substantially from better IO determinism.

# Shannon Direct-IO™ PCIe Flash G4i

**User Capacity**

**12.8TB**

**Access Latency**

**90/9us**

**IOPS**
Random 4KB read/write

**495/650K**

上海宝存信息科技有限公司
**Shannon Systems**