# Emerging Data-centric Storage Architectures

**Pankaj Mehra, Ph.D.**
**VP Storage Pathfinding**

**Samsung Semiconductor Inc.**

# Challenges of Data@Scale

## Bottlenecks

Processing power and
processing bandwidth

Metadata inefficiency of
object storage & retrieval

Wire protocol termination
for disaggregated flash

## Inefficiencies

Inability to deliver both
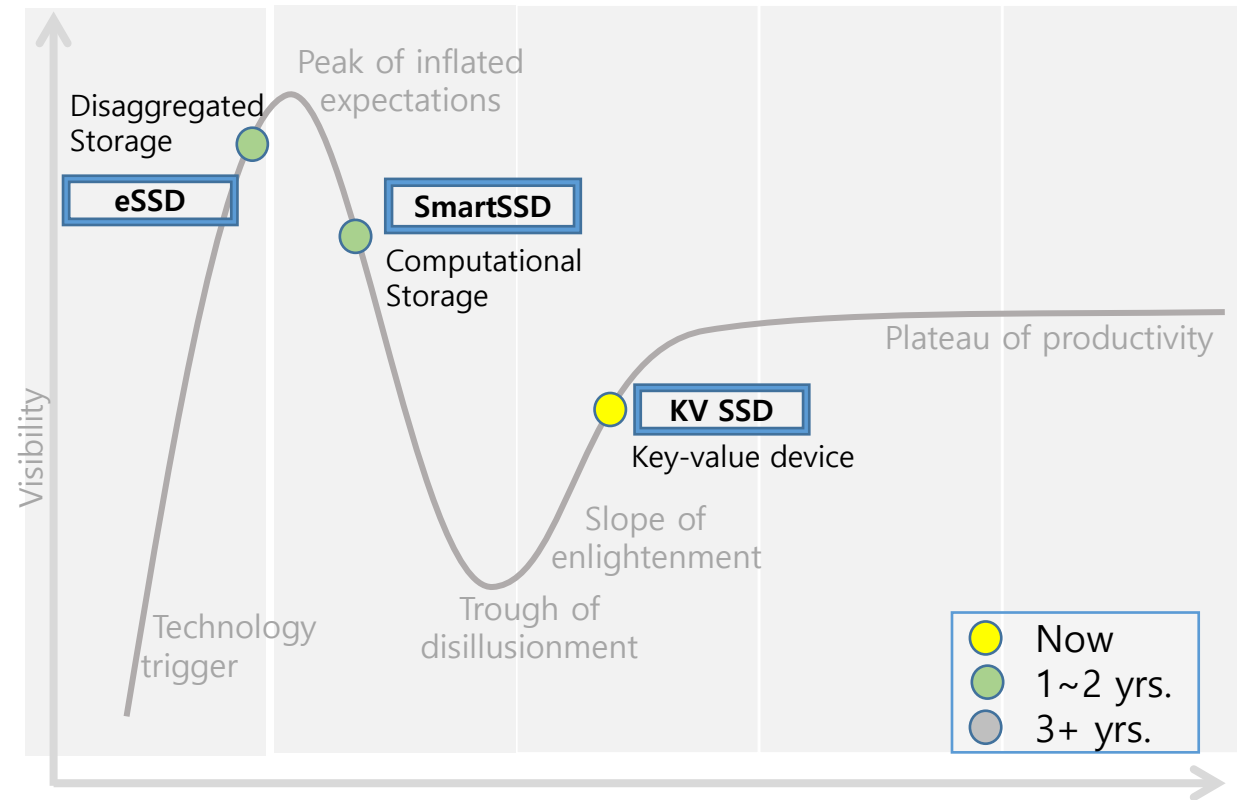performance and scale

Wasted endurance

Wasted memory BW

CPU overhead of I/O

CPU overhead of I/O
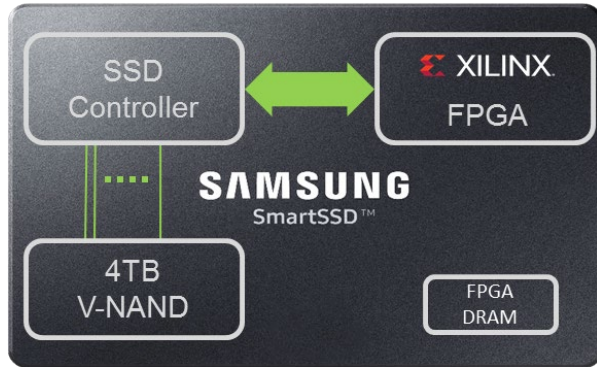virtualization

# Good Ideas, Already In-Play



- Virtualization offload
- SMRDB (since HDD days)
- DB filtering acceleration
- Storage NWconv (since FC)
- Active Disk (since HDD days)
- OSD (since HDD days)

*Why Revisit?*

Because in 2020, three distinct 25-y.o. ideas meet the SSD!

Disaggregated Storage

**eSSD**

Peak of inflated expectations

**SmartSSD**

Computational Storage

Plateau of productivity

**KV SSD**

Key-value device

Slope of enlightenment

Visibility

Trough of disillusionment

Technology trigger

Now
1~2 yrs.
3+ yrs.

# SmartSSD® CSD Scales to Accelerate Data-Rich Workloads

## SmartSSD U.2 Platform

SSD Controller ↔ XILINX FPGA

SAMSUNG SmartSSD™

4TB V-NAND

FPGA DRAM
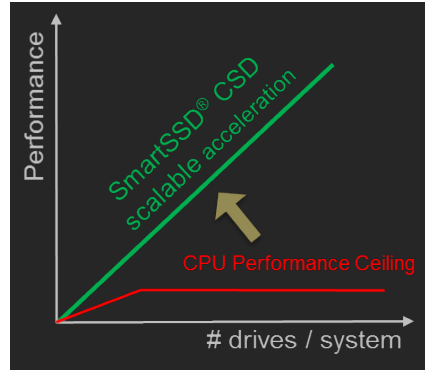
### Computational Storage

✓ **3 & 6 GBps internal BW per device:** Minimize external data movement

✓ **FPGA:** Each device has 3x~10x core equivalents for offload/acceleration

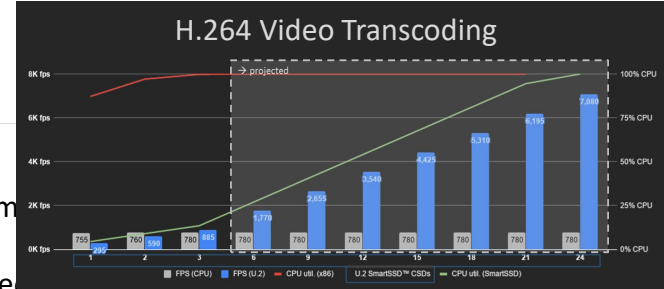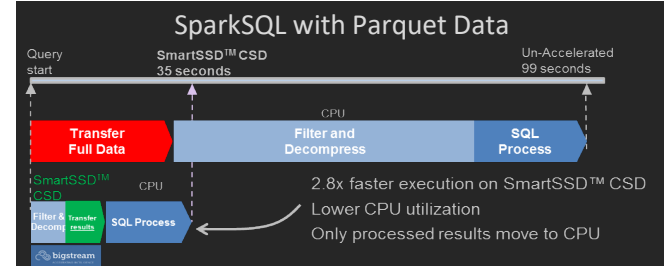✓ **4TB storage, 4 GB FPGA DRAM:** For Inline and Data@Rest processing

## Acceleration Concept

Performance vs # drives / system

SmartSSD® CSD scalable acceleration

CPU Performance Ceiling

### Scalable Performance

✓ **Near Data Processing:** Data format conversion, Filtering, Metadata management, DB Analytics, Video processing

✓ **New Services:** Secure content, Edge acceleration
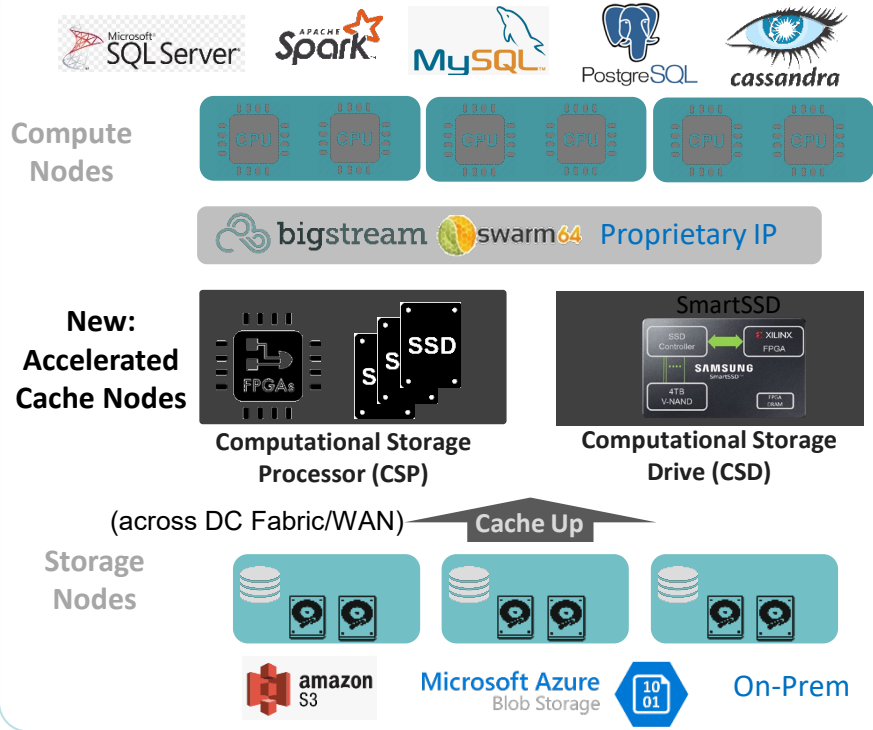
## Partner Solutions

### SparkSQL with Parquet Data

Query start — SmartSSD™ CSD 35 seconds — Un-Accelerated 99 seconds

Transfer Full Data | CPU — Filter and Decompress | SQL Process

SmartSSD™ CSD — Filter & Decomp | Transfer results | CPU — SQL Process | bigstream

2.8x faster execution on SmartSSD™ CSD
Lower CPU utilization
Only processed results move to CPU

### H.264 Video Transcoding

FPS (CPU) | FPS (U.2) | CPU util. (x86) | U.2 SmartSSD CSDs | CPU util. (SmartSSD)

### P2P Compression and Decompression

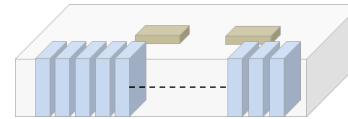| | CR | Throughput (MB/s) | CPU Eff. (MB/s/cpu) | PCIe Eff. | Mem Eff. |
|---|---|---|---|---|---|
| Compression: NoLoad-CSD | 2.85 | 2173 | 2173 | 1.72 | 1.42 |
| Compression: NoLoad-CSD w/ p2pdma | 2.85 | 2862 | 2938 | 1.01 | 0.06 |
| Decompression: NoLoad-CSD | 2.85 | 1823 | 2989 | 1.71 | 1.81 |
| Decompression: NoLoad-CSD w/ p2pdma | 2.85 | 2022 | 3690 | 1.01 | 0.14 |

# Computational Storage Use Cases Examples

- 3rd party and proprietary acceleration stacks run on Computational Storage to accelerate real-time analytics and regex searches for cybersecurity

## Analytics Cache Node

**Compute Nodes**

Microsoft SQL Server · Apache Spark · MySQL · PostgreSQL · cassandra

bigstream · swarm64 · Proprietary IP

**New: Accelerated Cache Nodes**

FPGAs · S S SSD

SmartSSD
SSD Controller · XILINX FPGA · SAMSUNG SmartSSD · 4TB V-NAND · FPGA DRAM

Computational Storage Processor (CSP)

Computational Storage Drive (CSD)

(across DC Fabric/WAN) ← **Cache Up**

**Storage Nodes**

amazon S3 · Microsoft Azure Blob Storage · 10 01 · On-Prem
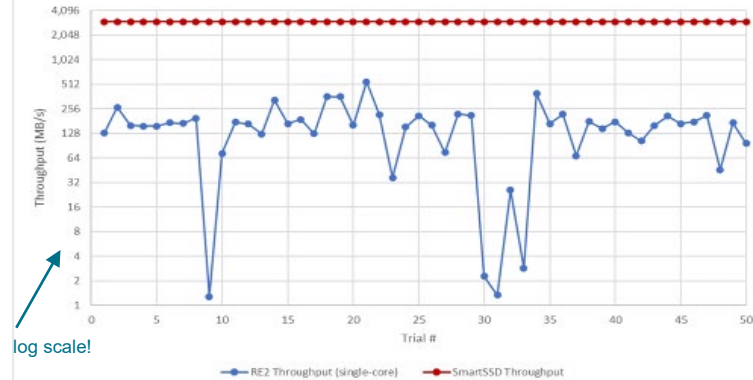
## RegEx Appliance

- **Throughput scales to large datasets and complex searches**
- **>10x throughput improvement compared to x86**

RegEx Appliance, 24x SmartSSD, 48TB
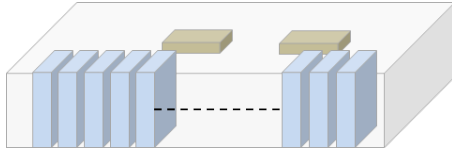
### RE2 vs SmartSSD Throughput
*(as function of expression complexity, set size, match density, near-match density, etc)*

Throughput (MB/s) — log scale!

Trial #

RE2 Throughput (single-core) · SmartSSD Throughput

# Samsung SmartSSD® Technology Roadmap

- Samples, development tools, partners solutions available for immediate PoC
- Customer PoC Test&Dev systems/support available from Samsung and partners
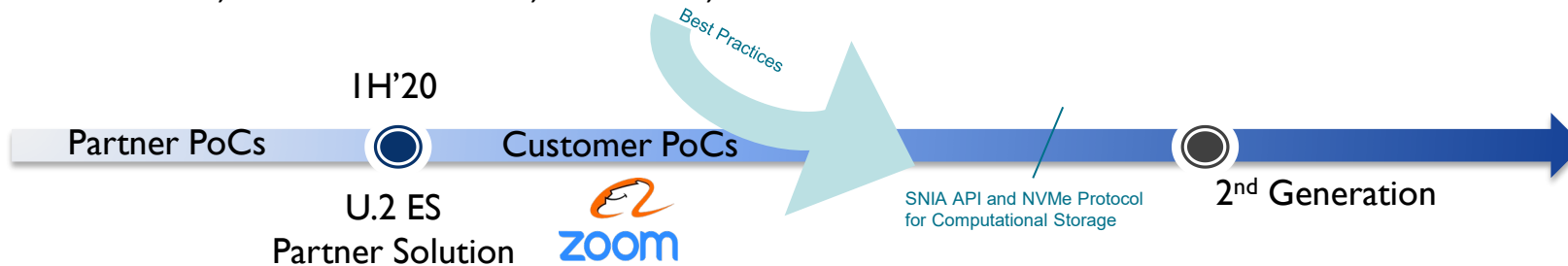


**v1.0 SmartSSD® U.2 CSD**

**U.2 FF: Scale Processing to 24 ~ 48 devices**
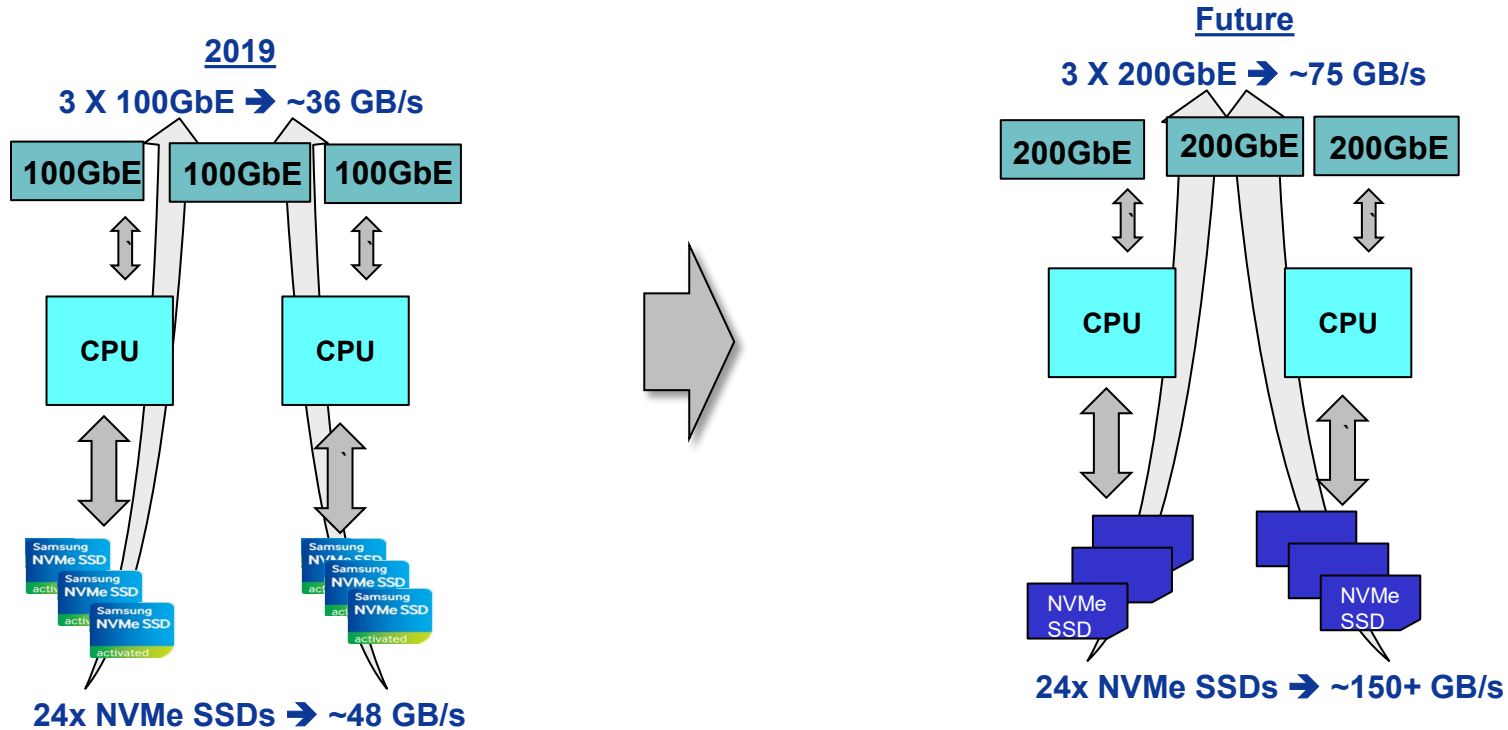4TB, PCIe Gen3x4 External, ~530K LUTs,

**Next Gen SmartSSD® CSD**

Customers requirements: Integration, Interfaces, FF, workloads

Best Practices

1H'20

Partner PoCs          Customer PoCs

U.2 ES
Partner Solution          zoom

SNIA API and NVMe Protocol
for Computational Storage          2nd Generation

# Ethernet SSD targets IO bottleneck in Storage Chasses



**2019**

**3 X 100GbE ➔ ~36 GB/s**

100GbE    100GbE    100GbE

CPU        CPU

Samsung NVMe SSD
activated

**24x NVMe SSDs ➔ ~48 GB/s**

**Future**

**3 X 200GbE ➔ ~75 GB/s**

200GbE    200GbE    200GbE

CPU        CPU

NVMe SSD    NVMe SSD

**24x NVMe SSDs ➔ ~150+ GB/s**

**CPU and IO bottleneck for storage throughput performance**

# NVMe-oF SSD based EBOF

## Conventional NVMe JBOF

**Storage Head Nodes or Application Servers**

**Ethernet Switch**

PCIe Switch

PCIe

| NVMe SSD | NVMe SSD | NVMe SSD | NVMe SSD |

**NVMe JBOF**

❑ **Pros**
- ✓ Enables disaggregation of NVMe SSDs
- ✓ Management & Storage Services
- ✓ Utilizing existing storage & server architectures

❑ **Cons**
- ➢ Non-scalable Storage Controller - PCIe single root constraint
- ➢ Bandwidth Limitation
  - • CPU, PCIe, NW Constraints
- ➢ Power and Thermals

## NVMe-oF EBOF

| Storage Controllers | App. Server Web | App. Server Database | App. Server Analytics |

**Hypervisor, Container**

**Ethernet Switch**

Ethernet Switch | Ethernet Switch

**Ethernet**

| NVMe-oF SSD | NVMe-oF SSD | • • • | NVMe-oF SSD | NVMe-oF SSD |

**NVMe-oF EBOF**

❑ **Pros**
- ✓ High Bandwidth
- ✓ Scaled Linearly (Ethernet)
- ✓ Sharable via NVMe-oF
- ✓ Less power
- ✓ Lower latency

❑ **Cons**
- ➢ New platform architecture
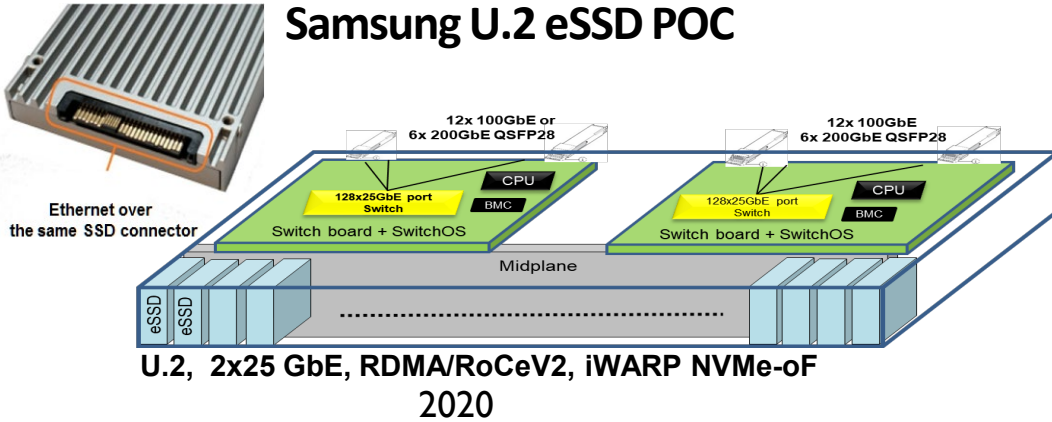- ➢ Management of Storage Services & Network Devices

**NVMe-oF EBOF can address bandwidth, scalability, and flexibility**

# Samsung Ethernet SSD Technology Roadmap

- Samples, development tools, partners solutions available for immediate PoC
- Customer PoC Test&Dev systems/support available from Samsung and partners

**Samsung U.2 eSSD POC**

**Next Gen Ethernet SSD**

Ethernet over
the same SSD connector

12x 100GbE or
6x 200GbE QSFP28

12x 100GbE
6x 200GbE QSFP28

CPU

128x25GbE port
Switch

BMC

Switch board + SwitchOS

CPU

128x25GbE port
Switch

BMC

Switch board + SwitchOS

Midplane

eSSD  eSSD

**U.2,  2x25 GbE, RDMA/RoCeV2, iWARP NVMe-oF**

**2020**

SAMSUNG

Customers requirements:
NVMe-oF/TCP, HW Offloads
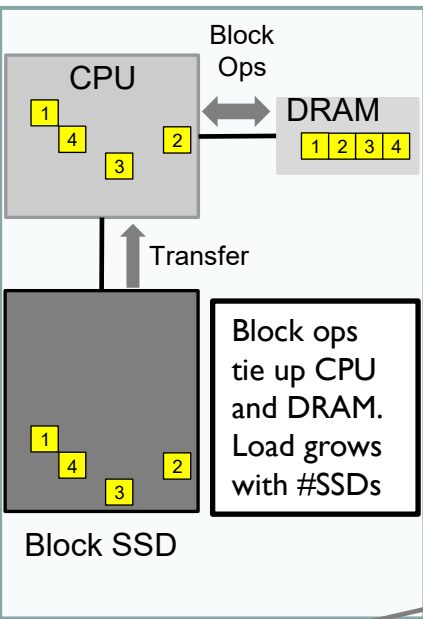FF, Congestion Control (ECN/PFC),
Security, Management

PoCs

U.2 ES
POC
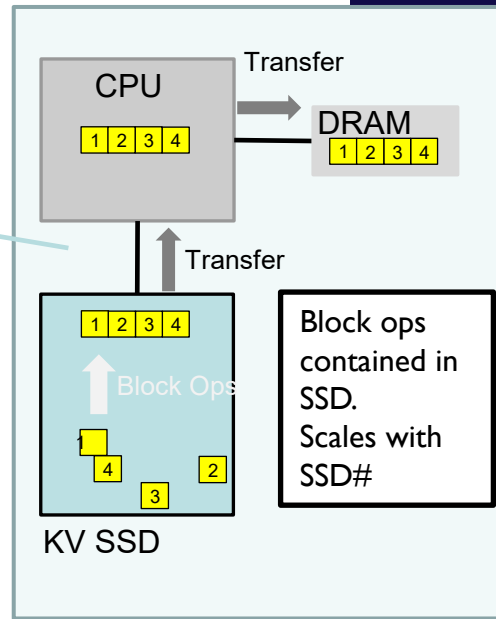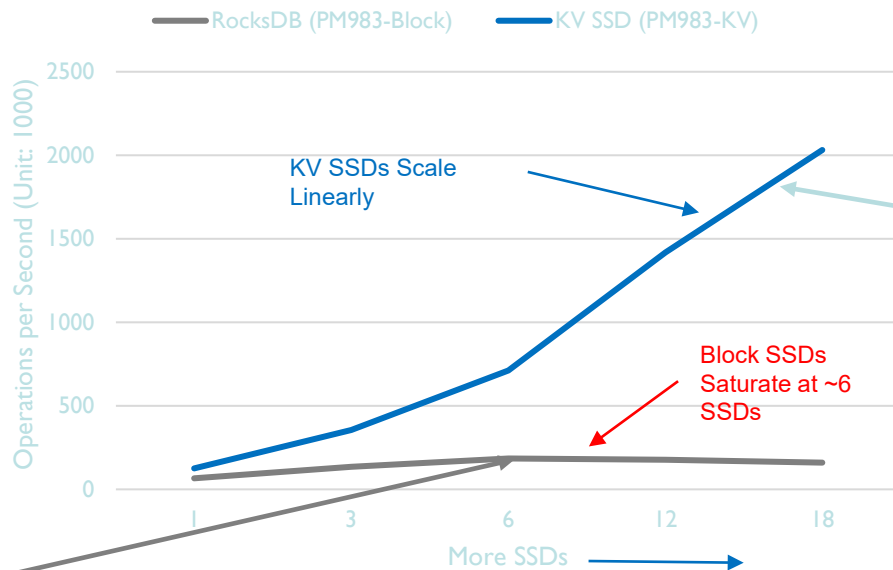
2nd Generation

# KV SSD is about Efficient Block Operations

Block operations on CPUs ⇨ Bottlenecks, Scaling Inefficiency.
KV SSD offloads Block operations from CPUs.
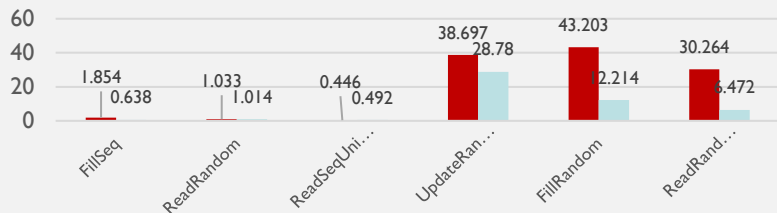


Server Scalability with Increasing # of SSDs

Block Ops

CPU

DRAM

Transfer

Block ops tie up CPU and DRAM. Load grows with #SSDs

Block SSD

Block SSDs

RocksDB (PM983-Block)
KV SSD (PM983-KV)

KV SSDs Scale Linearly

Block SSDs Saturate at ~6 SSDs

More SSDs

CPU

Transfer

DRAM

Transfer

Block Ops

KV SSD

Block ops contained in SSD.
Scales with SSD#

# KV-optimized SW shows Multiple Efficiency Benefits

## Average Latency (us)

■ RocksDB

| | |
|---|---|
| FillSeq | 1.854 / 0.638 |
| ReadRandom | 1.033 / 1.014 |
| ReadSeqUni... | 0.446 / 0.492 |
| UpdateRan... | 38.697 / 28.78 |
| FillRandom | 43.203 / 12.214 |
| ReadRand... | 30.264 / 6.472 |

Average Latency is up to 60% better on KVRocks

## Insertion operations/sec

KVRocks   RocksDB

400, 149194.76

400, 51199.15

Keys Count (Millions)

Insertion Operations / sec

## Application WAF

■ RocksDB  ■ KVRocks

Overwrite
FillRandom
FillSeq

0    5    10    15    20    25

## SSD Lifetime

■ RocksDB  ■ KVRocks

Overwrite
FillRandom
FillSeq

0   2   4   6   8   10   12   14

SSD lifetime is more than 11x for KVRocks compared to RocksDB ⇨ **11x less CapEx**

# Samsung KV SSD Technology Roadmap

- Stack open sourced at https://github.com/OpenMPDK/KVSSD

**U.2 KV SSD**

**U.2 FF: Scale Processing to 24 devices**

4TB, PCIe Gen3x4 External

**Next Gen KV SSD**

Your requirements?
(Integration, Interfaces, FF, Workloads)

1H'20

Software dev

Partner Testing
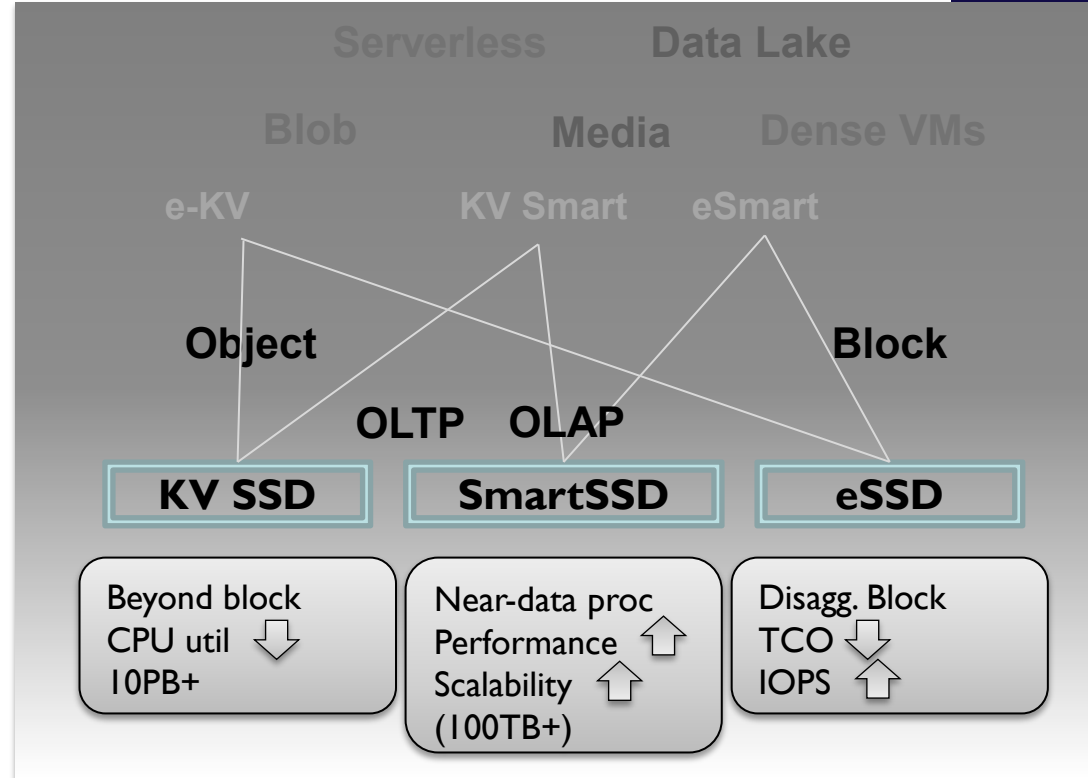
Rocks DB repl.
KVCeph
Minio

Partner Solution

2nd Generation

# Benefits Summary

| | SmartSSD | Ethernet SSD | Key-Value SSD | Zoned Name Spaces |
|---|:---:|:---:|:---:|:---:|
| **Application Awareness** | ✓ | | ✓ | ✓ |
| **Acceleration** | ✓ | | | |
| **Reduce data-related CPU load** | ✓ | | ✓ | |
| **Improved Write Endurance** | | | ✓ | ✓ |
| **Fewer protocol terminations** | ✓ | ✓ | | |
| **Min device virtualization o/h** | | | | ✓ |
| **Fewer stack translations** | | ✓ | ✓ | |
| **Metadata Optimization** | | | ✓ | |
| **Scaling Data Bandwidth** | ✓ | | | |
| **Saving L2-to-Memory BW** | ✓ | | ✓ | |
| **Control@Scale (IODT, QoS)** | | | | ✓ |
| **Maximize #SSDs/chassis** | ✓ | ✓ | ✓ | ✓ |

# Possible Convergence

| Host Interface | Addressing | Accelerator |
|---|---|---|
| PCIe | Block | None |
| Ethernet | ZNS | FPGA |
| | Key-Value | |

Serverless     Data Lake

Blob     Media     Dense VMs

e-KV     KV Smart     eSmart

**Object**     **Block**

OLTP   OLAP

**KV SSD**     **SmartSSD**     **eSSD**

Beyond block
CPU util ⬇
10PB+

Near-data proc
Performance ⬆
Scalability ⬆
(100TB+)

Disagg. Block
TCO ⬇
IOPS ⬆

**Please take a moment
to rate this session.**

**Your feedback matters to us.**

- Bullet one
  - Bullet two
    - Bullet 3
      - Bullet 4

# Dr. Pankaj Mehra

Vice President,
Storage Pathfinding

Samsung Semiconductor Inc.

- Subhead
  - Example 1
  - Example 2
- Subhead