# Storage in the New Age of AI/ML

**Young Paik**

**Sr Director Product Planning**

**Samsung**

May 21, 2019

COLLABORATE. INNOVATE. GROW.

SAMSUNG

# Legal Disclaimer

This presentation is intended to provide information concerning SSD and memory industry. We do our best to make sure that information presented is accurate and fully up-to-date. However, the presentation may be subject to technical inaccuracies, information that is not up-to-date or typographical errors. As a consequence, Samsung does not in any way guarantee the accuracy or completeness of information provided on this presentation.

The information in this presentation or accompanying oral statements may include forward-looking statements. These forward-looking statements include all matters that are not historical facts, statements regarding the Samsung Electronics' intentions, beliefs or current expectations concerning, among other things, market prospects, growth, strategies, and the industry in which Samsung operates. By their nature, forward-looking statements involve risks and uncertainties, because they relate to events and depend on circumstances that may or may not occur in the future. Samsung cautions you that forward looking statements are not guarantees of future performance and that the actual developments of Samsung, the market, or industry in which Samsung operates may differ materially from those made or suggested by the forward-looking statements contained in this presentation or in the accompanying oral statements. In addition, even if the information contained herein or the oral statements are shown to be accurate, those developments may not be indicative developments in future periods.

SAMSUNG

# Speaker Disclaimer

Sometimes accuracy is the enemy of the truth

# AI/ML Workflow – So Simple

Logs

DBs

Real-time streams

Images

Video

Audio

IoT

Genetic Data
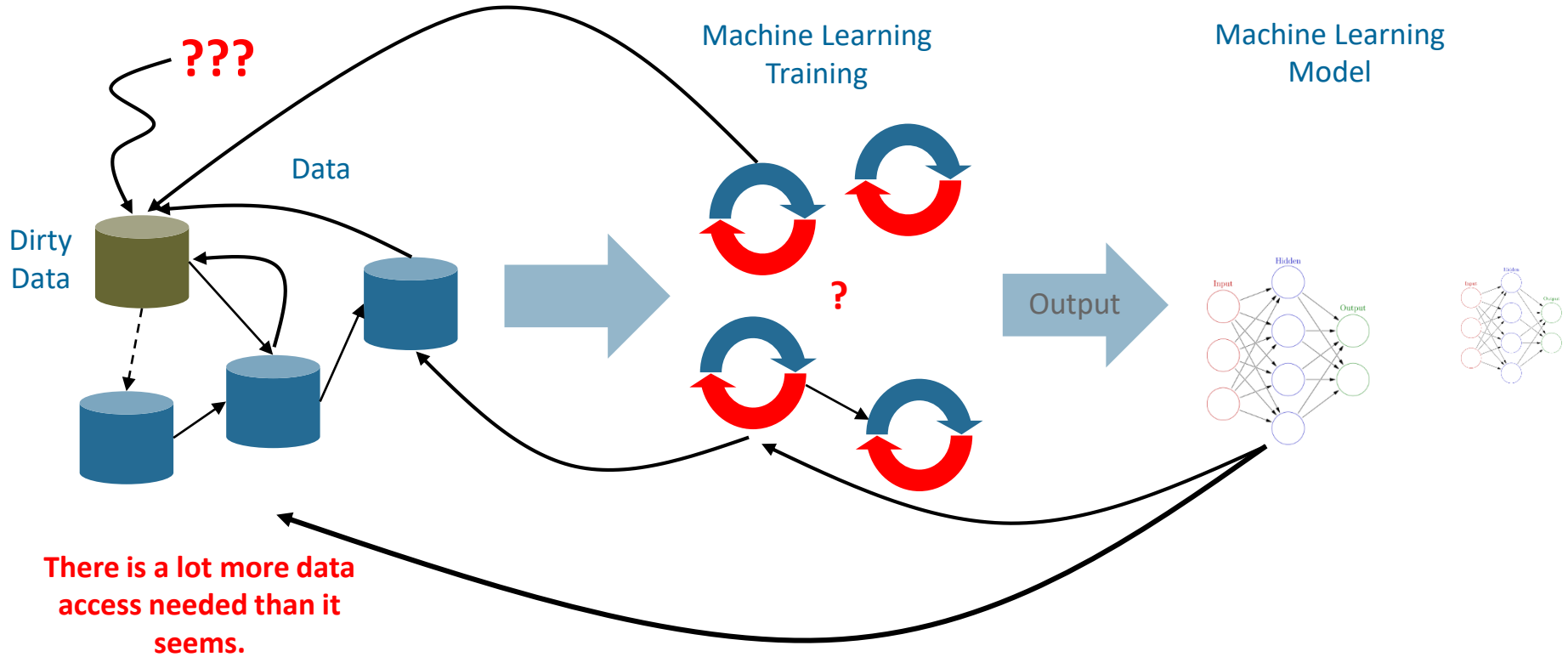
Data

Machine Learning Training

Output

Machine Learning Model

**But is it really this simple?**
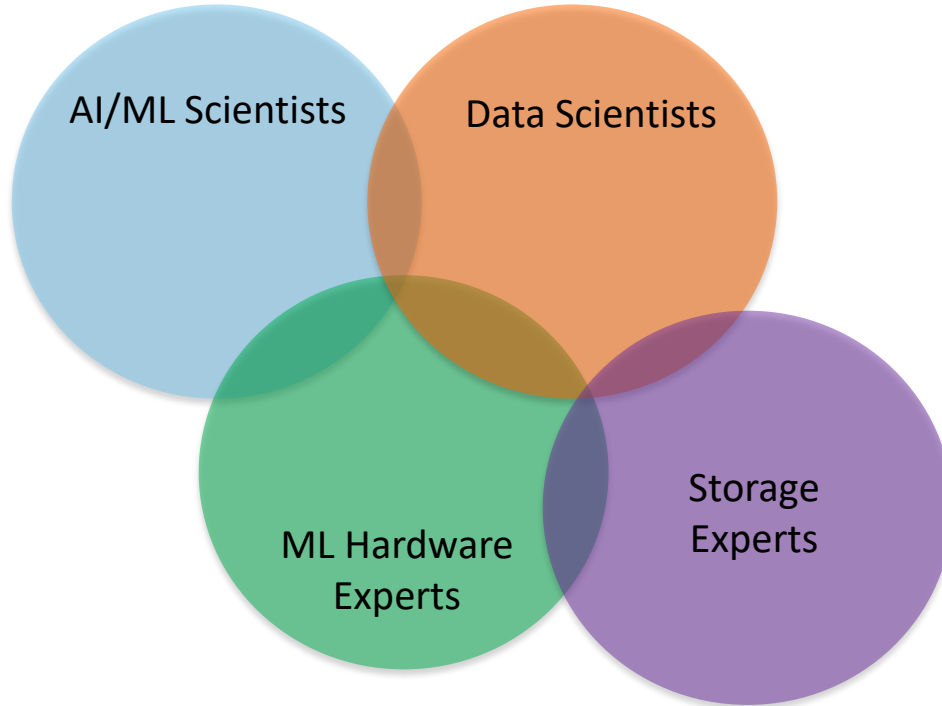
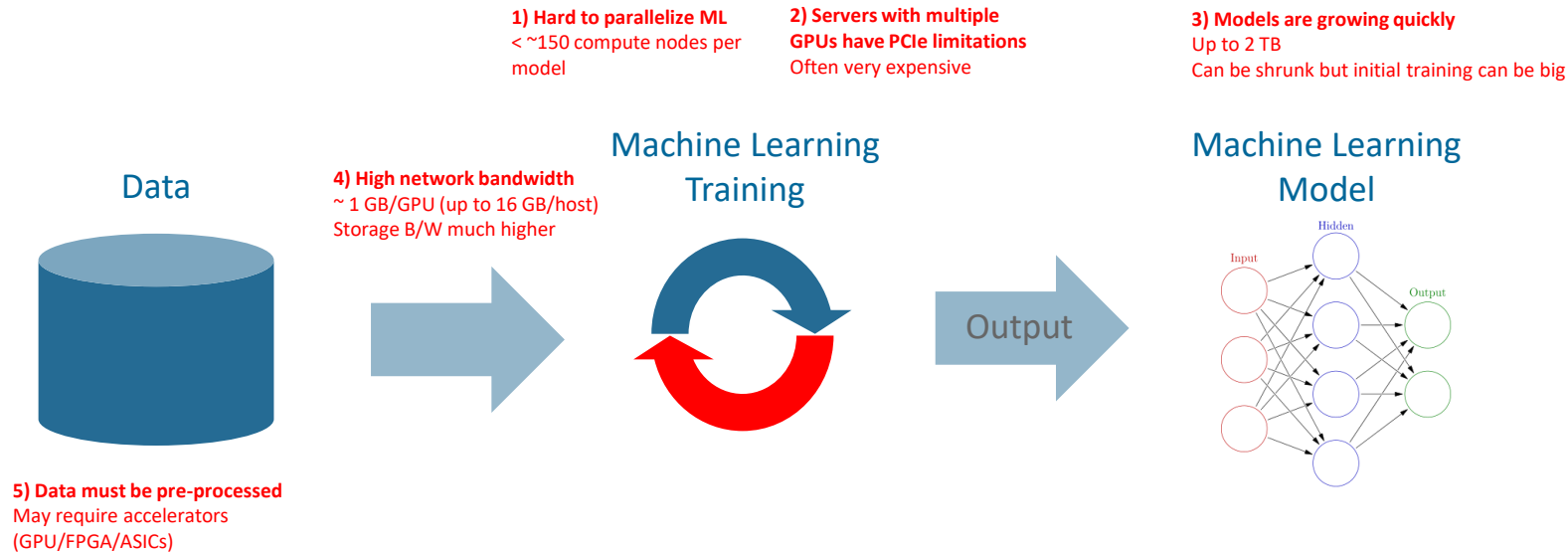# AI/ML Workflow – It's Never Easy

# Disparate Groups of Experts

Skill sets are highly specialized, often without overlapping skill sets

SAMSUNG

# Artificial Intelligence Workflow – Major Challenges

1) Hard to parallelize ML
< ~150 compute nodes per model

2) Servers with multiple GPUs have PCIe limitations
Often very expensive

3) Models are growing quickly
Up to 2 TB
Can be shrunk but initial training can be big

Data

4) High network bandwidth
~ 1 GB/GPU (up to 16 GB/host)
Storage B/W much higher

Machine Learning Training

Output

Machine Learning Model

5) Data must be pre-processed
May require accelerators
(GPU/FPGA/ASICs)

**What does preprocessing look like?**

# Artificial Intelligence Workflow – Facial Recognition

**Images**

**Facial Recognition Training**

**Facial Recognition Model**

Output

Deep Learning models need the same facial form

AI/ML Training servers may cost up to $400K

• Photo by rawpixel.com from Pexels

SAMSUNG

# Facial Recognition Example of Preprocessing



| 1) Find faces | 2) Extract faces | 3) Resize image and color | 4) Rotate face | 5) Extract features |
|---|---|---|---|---|

(My sincere apologies to the model for this rendering)

Photo by rawpixel.com from Pexels

To recognize the identity of a face, you must first isolate every face.

Training must work on individual faces

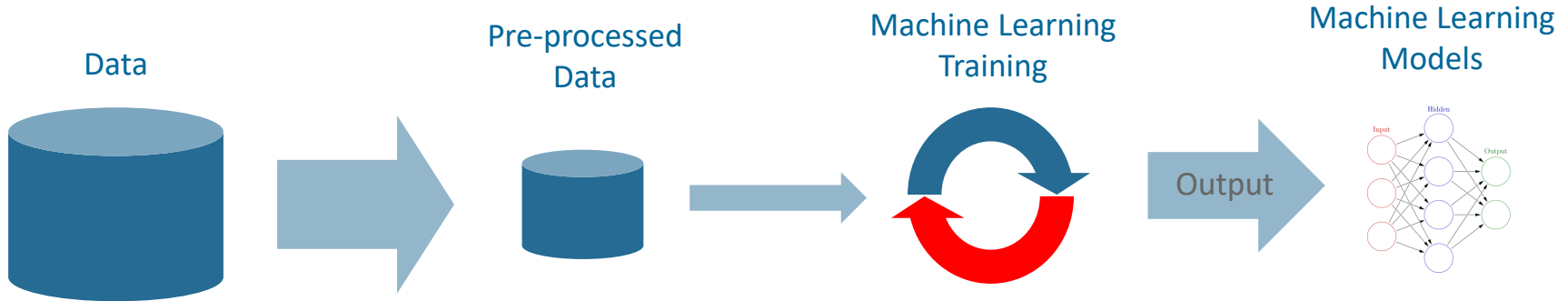Images must conform to the same pixel and color resolution

Face must be front (there are algorithms that do this)

You can now extract the facial features and begin the training

## All of this is parallelizable and does not need to be done on the training server

# Artificial Intelligence Workflow – Add Preprocessing

Data

Pre-processed Data

Machine Learning Training

Machine Learning Models

Output
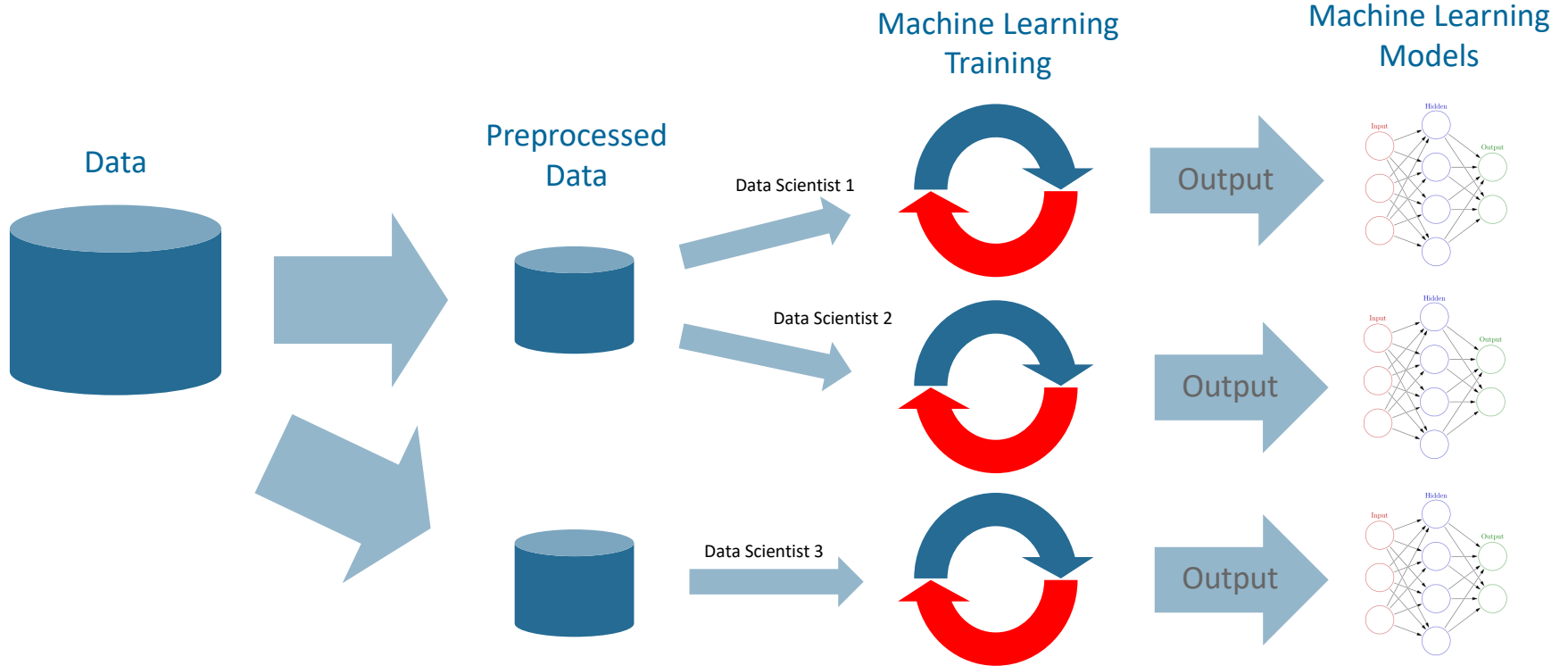
**More Complicated Issues**

Multiple AI Scientists

Improved data processing

Dealing with long training times

SAMSUNG
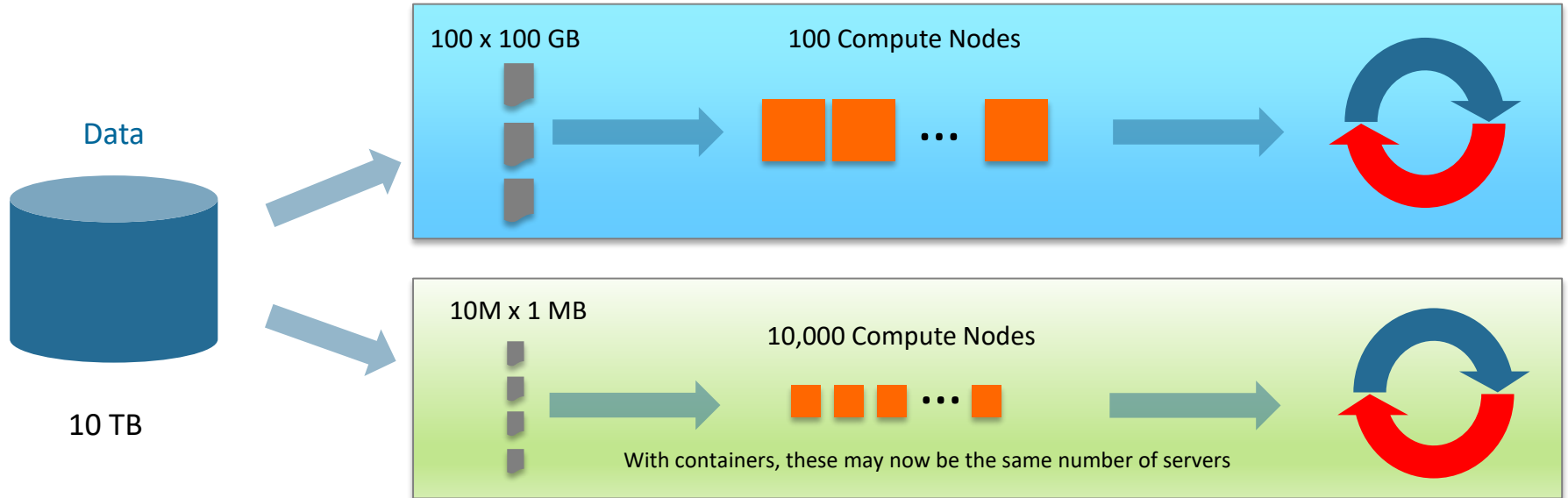
# Multiple Data Scientists

Data scientist 1 and 2 want the same features, but different models
Data scientist 3 is trying a new experiment and must start from raw data

# Dealing With Long Training Times
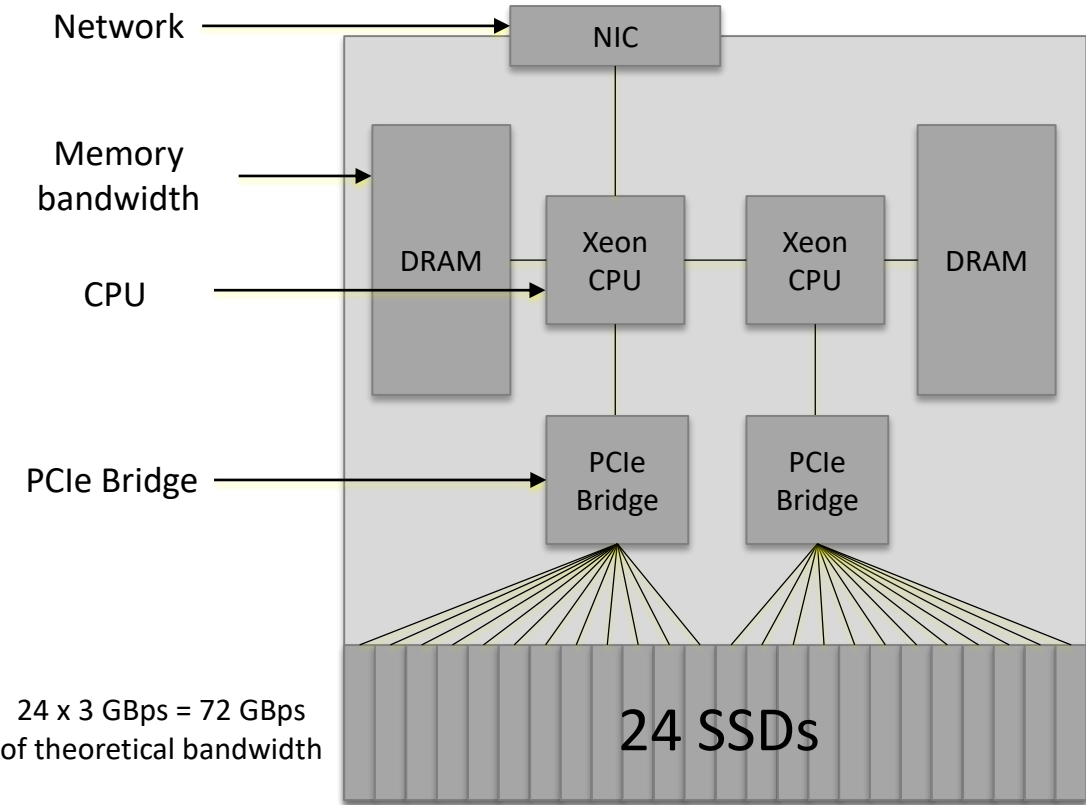
Training times may take weeks.
How can we deal with changes in workload dictated by changing priority?



**10 TB**

Challenges:
- Minimum size for jobs (not all jobs can be shrunk)
- Scheduling is huge → Kubernetes
- Jobs are not always parallelizable (database joins)

# Data Flow Limits of Modern Storage

Network

NIC

Memory
bandwidth

DRAM

Xeon
CPU

Xeon
CPU

DRAM

CPU

PCIe Bridge

PCIe
Bridge

PCIe
Bridge

24 SSDs

24 x 3 GBps = 72 GBps
of theoretical bandwidth

Modern SSDs are limited by
server architecture

Samsung has looked into 2
different technologies:
- KV SSD
- SmartSSD
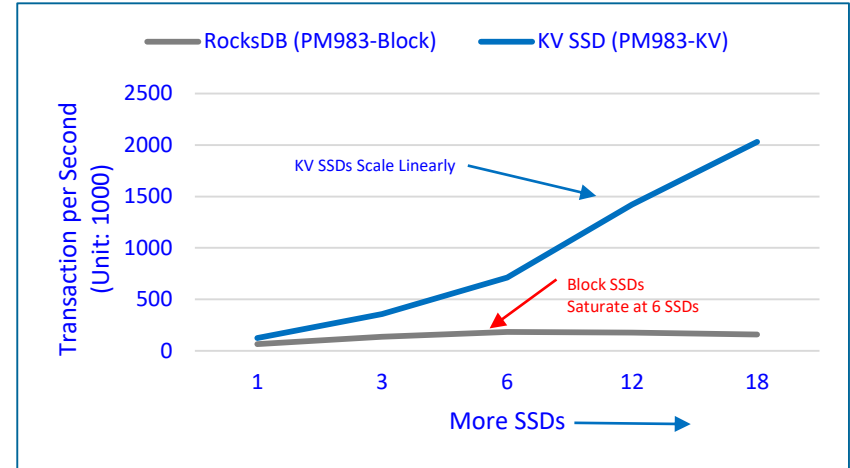
# KV SSD - Motivation

## KV API now a SNIA Specification

https://www.snia.org/tech_activities/standards/curr_standards/kvsapi

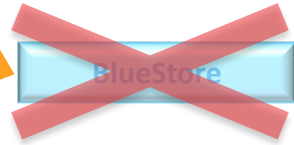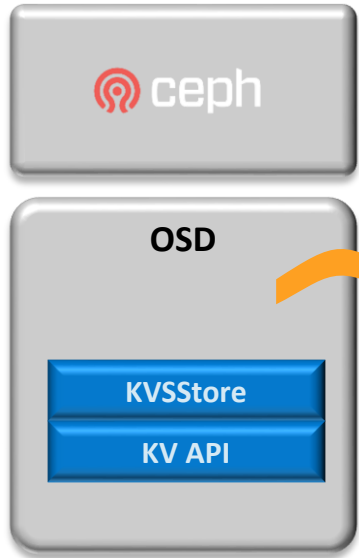| | Block SSD | KV SSD |
|---|---|---|
| CPU | Overloaded with block and compaction | Freed for other tasks |
| Scalability | Limited to 4-6 SSDs/host | Linear performance with 18+ SSDs/host |
| Disk utilization | Must leave room for compaction | GC managed internally |
| SSD Lifetime | High WAF | Low WAF leads to greatly improved SSD lifetime |



* Testing was done on a server with 2 x Intel Xeon E5-2600 v5 servers with 384 GB of DRAM, and 18 PM983 (in block or KV mode) SSDs
** Workload: 4KB uniform random writes

Main Use Cases:
- Object storage
- NoSQL databases

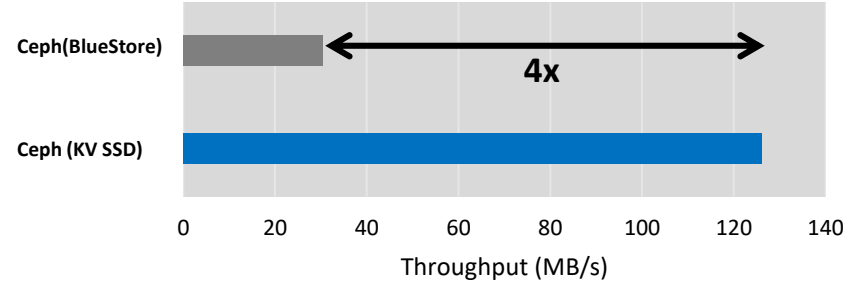SAMSUNG

# KV SSD – Direct Use on Ceph



**OSD**

**KVSStore**

**KV API**

KVSStore uses the newly Open Sourced KV API to access the KV SSD

*https://github.com/OpenMPDK/KVSSD

KV

SAMSUNG
Solid State Drive

Biggest challenge is that this requires a change in software.

## Higher Throughput



Ceph(BlueStore)

**4x**

Ceph (KV SSD)

Throughput (MB/s)

0  20  40  60  80  100  120  140
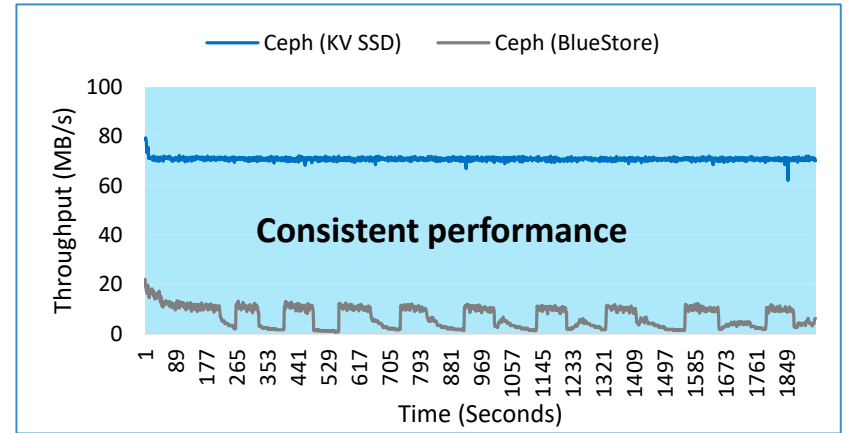
* 4096 block write Default (Sharded), 8 clients 2 OSDs- queue depth 128
* Testing was done on two servers with 2 x Intel Xeon E5-2695 v4 CPUs with 128 GB of DRAM, and a PM983 (in block or KV mode) SSD with 40 GbE



Ceph (KV SSD)    Ceph (BlueStore)

Throughput (MB/s)

100
80
60
40
20
0

**Consistent performance**

Time (Seconds)

1  89  177  265  353  441  529  617  705  793  881  969  1057  1145  1233  1321  1409  1497  1585  1673  1761  1849

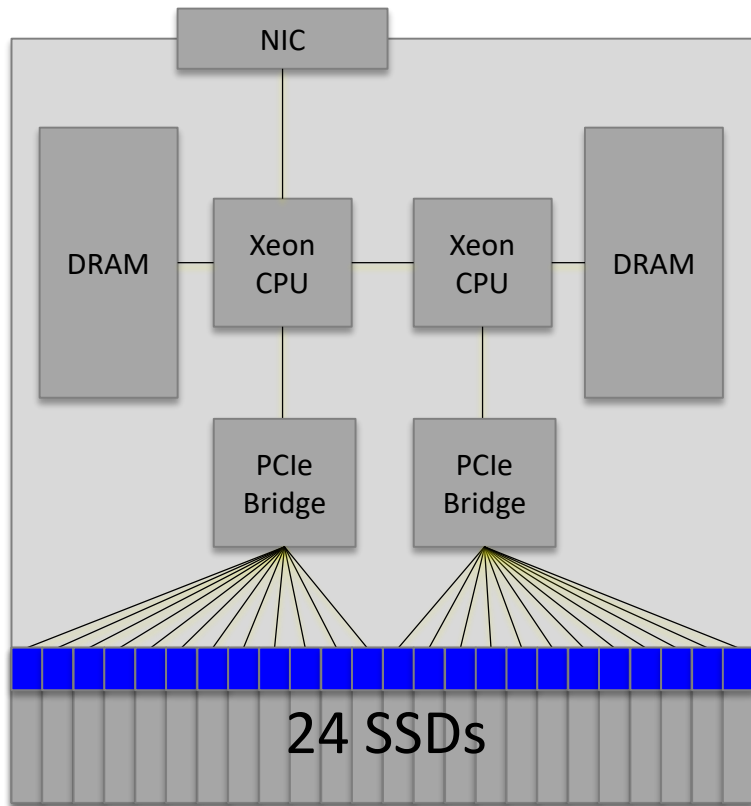* 4096 block write Default (Sharded), 1 client 1 OSD - queue depth 128
* Testing was done on a server with 2 x Intel Xeon E5-2695 v4 CPUs with 128 GB of DRAM, and a PM983 (in block or KV mode) SSD with 40 GbE

# SmartSSD-based Server Architecture



SmartSSDs process
data in-storage

Allows:
- Pre-filtering
- On-disk transcoding
- Compression
- ...

Challenges:
- Encryption
- RAID/Erasure Coding
- New programming model

Compute occurs on storage
Parallel scans at full speed of SSDs
CPUs freed for additional work

NIC

DRAM

Xeon CPU

Xeon CPU

DRAM

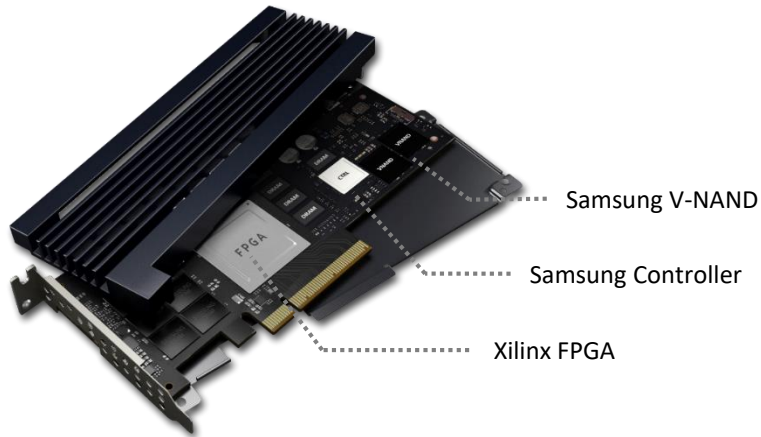PCIe Bridge

PCIe Bridge

24 SSDs

SAMSUNG

# SmartSSD

## SmartSSD PM983F announced at Samsung Tech Day 2018

### PM983F AIC
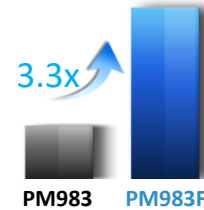
- SmartSSD PCIe add-in card
- Shown successfully integrated with Bigstream
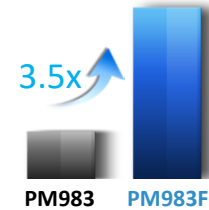- Several data-intensive workloads easily ported



Samsung V-NAND

Samsung Controller

Xilinx FPGA

### PoC Results

- For I/O-bound workloads, SmartSSD showed 3x to 4x better performance with scalability

Financial BI (VWAP[1])
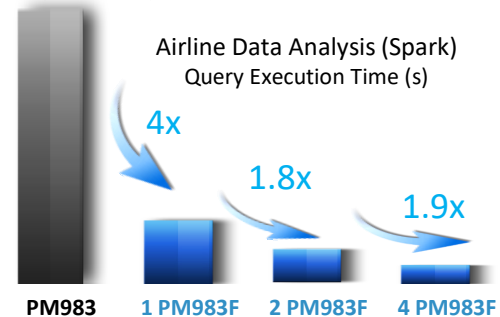Throughput (MOPS)

3.3x

PM983    PM983F

Database (MariaDB)
TPC-H Score, Geo.Mean

3.5x

PM983    PM983F

* VWAP: Volume Weighted Average Price

Airline Data Analysis (Spark)
Query Execution Time (s)

4x

1.8x

1.9x

PM983    1 PM983F    2 PM983F    4 PM983F

SAMSUNG

COLLABORATE. INNOVATE. GROW.

# New Technologies Not Covered

| Technology | Description | Pros | Cons |
|---|---|---|---|
| Nvidia GPUDirect | GPUs can directly access another PCIe device | Bypasses CPU and system memory | Some people use system memory as a cache |
| NVMe over Fabric | Allows for very low latency to network-attached storage with RDMA latencies | Gives performance similar to direct-attach | Requires very solid network coordination |
| SmartNICs | These NICs have CPU offload facilities. Many have the ability to handle Reed-Solomon. | Low latency at a much lower price point. | Still very new |

# Thank You

Young.Paik@Samsung.com

SAMSUNG

COLLABORATE. INNOVATE. GROW.

# How Important Is It To Fix Dirty Data?

**Scenario:** One healthcare insurance company was looking at data on charges for treatments.

We tested by looking at diseases by code and tried to guess what the disease was.

| Disease Code 1 | |
|---|---|
| Average age: | 63 |
| Gender | |
| Male | 33% |
| Female | 66% |
| Unspecified | 1% |
| | |
| Diagnosis: Osteoperosis | |

| Disease Code 2 | |
|---|---|
| Average age: | 47 |
| Gender | |
| Male | 98% |
| Female | 0.5% |
| Unspecified | 1.5% |
| | |
| Diagnosis: ??? | |

Moral to the story: **It is important to thoroughly process data.**

**This requires much more storage I/O than people think.**

# MINIO + KV SSD Object Storage Performance

**DFSIO Benchmark**

**8 x Spark Node Cluster**

Dell 740xd
Intel 6152 (2.1GHz)
384 GB DDR4
1 x 100 GbE

| Spark Node | Spark Node | Spark Node | Spark Node | Spark Node | Spark Node | Spark Node | Spark Node |

**100GbE** — **Network SW** — **S3 API Protocol**

**MINIO** NKV API | **MINIO** NKV API | **MINIO** NKV API | **MINIO** NKV API

KV   PM983 KV

**4 x MINIO + KV SSD Cluster**

2 x Intel 6152 (2.1 GHz)
12 x 4 TB KV SSD
1 x 100 GbE
384 GB DDR4 (2400 MHz)
12 + 4 Erasure Code

## Minio Bandwidth on DFSIO on Spark with 4 Nodes



Bandwidth (GB/s)

100 MB: RD 27.44, WR 6.26
1000 MB: RD 24.77, WR 8.35

File Size

\* Performance tests were run with cache enabled for directory listing

SAMSUNG