# CINABRO:
## a Software Driven,
## Open Flash Array Architecture for
## Scalable Cloud Storage Services

Sungjoon Ahn, VP of Engineering, Circuit Blvd., Inc.
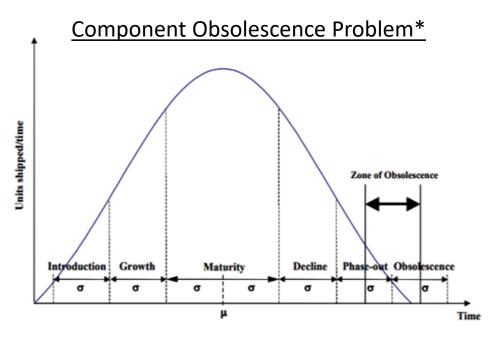
OPEN. FOR BUSINESS.

OCP
SUMMIT

# Motivation

## 1. Serves diverse cloud storage requirements

- Data center workloads are dynamic, diverse and constantly evolving
- Data center SSDs typically run 3 to 5 years after rigorous qualification process
- SSD FW update is expensive and usually limited to critical bug fixing

## 2. Streamlines flash memory deployments

- SSD designs optimized for single self contained units
- Data center SSDs often have old generation NANDs
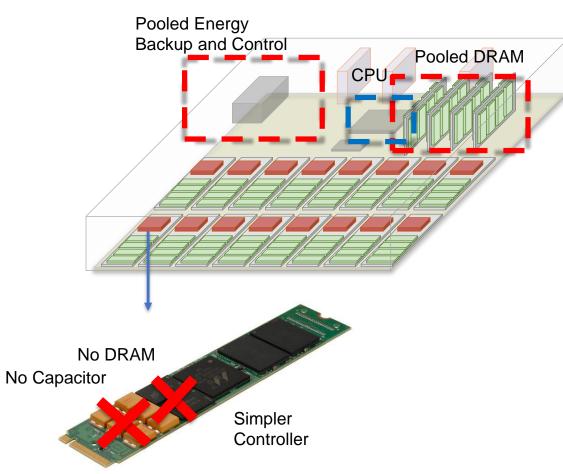- Need for deploying latest NANDs in scale



Component Obsolescence Problem*

# Solution

- Lower Total BOM and Simpler SSD Device Architecture
- Large portions of SSD intelligence run on host server CPU



Pooled Energy Backup and Control

Pooled DRAM

CPU

No DRAM
No Capacitor

Simpler Controller

| Category | Conventional SSD | Cinabro SSD |
|---|---|---|
| 1GB DRAM | $$ | X |
| 1TB NAND | $$$ | $$$ |
| SoC Controller | $$ | $ |
| Capacitors (x20) – Power Loss Protection | $ | X |
| Power Consumption | High | Low |
| Development Complexity | High | Low |

3

CIRCUITBLVD
Change your memories. Change your life.

# Cinabro™ System Architecture

## Disaggregated and composable All Flash Array based on COTS server



NVMe-over-Fabrics(RDMA)

to/from ToR (Eth/IB)

Out-of-Band Management

NIC/SmartNIC

Pooled Energy Backup and Control

Pooled DRAM

CBOS™

CPU

PCIe Switch

Open-Channel SSD

CBBridge™

NAND Modules

CIRCUITBLVD
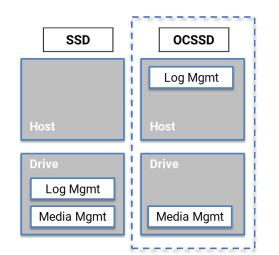Change your memories. Change your life.

# NAND Flash Interface

## Leverage OCSSD Standard to Provide Optimized Solution with Simpler ASIC



Physical flash exposed to the host (Read, write, erase)

Host
- Data placement
- IO Scheduling
- Over-provisioning
- Garbage collection
- Wear levelling

SSD / OCSSD

Host

Log Mgmt

Drive

Log Mgmt

Media Mgmt

- Open-Channel SSD (OCSSD)
  - Standard NVMe based protocol
  - Facilitates host FTLs and good fit for cloud providers

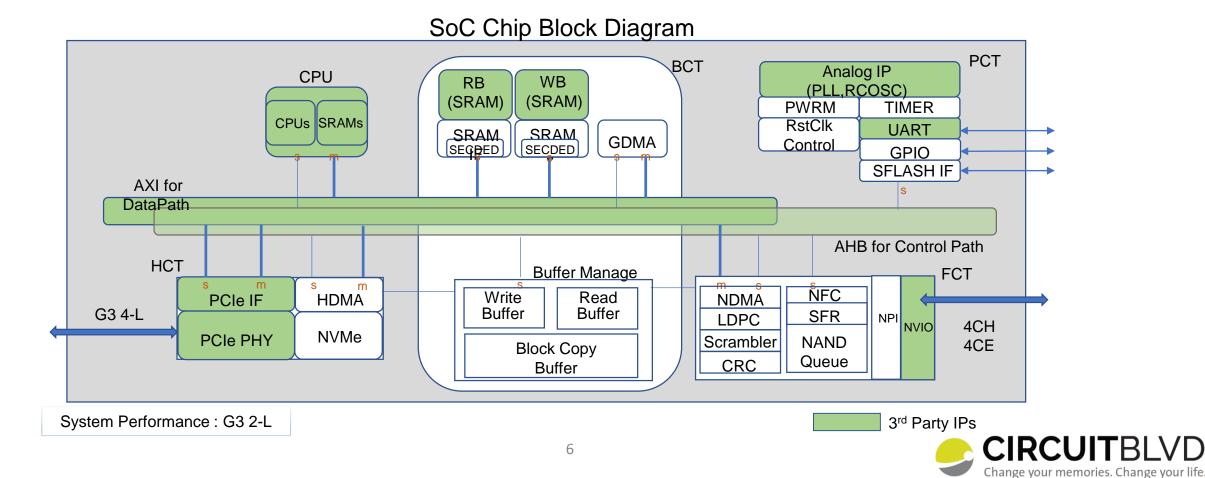- Optimized protocol translation between host and NAND interface

- Performance acceleration and reliability enhancement features for 3D NAND TLC/QLC

- Cost and power efficient ASIC design

CIRCUITBLVD
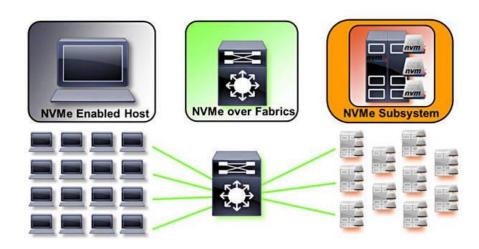Change your memories. Change your life.

# CBBridge™ OCSSD Controller

## Simple but robust SoC handling essential NAND media functions

- Open-Channel SSD spec and additional features for cross layer optimizations
- 28nm process technology accommodates 96+ layer 3D toggle3 TLC/QLC NAND with LDPC

### SoC Chip Block Diagram

# Network Interface

Leverage New Standard for Networked Storage Interfaces

- NVMe-over-Fabrics (NVMe-oF)
  - Faster access between hosts and storage systems
  - Much lower latency than iSCSI

- Flexible system design to support various fabrics of NVMe-oF standard (Ethernet, Infiniband, etc.)

- Open architecture allows incorporating new system technologies (e.g. SmartNIC, FPGA acceleration, SDN)

- Seamless integration with Open-Channel SSDs

**CIRCUIT**BLVD
Change your memories. Change your life.

# Software Design

## Advanced Open Source Software Optimized for All Flash Array

Linux LightNVM
SPDK/DPDK
RocksDB
OpenStack
Ceph
Docker
Kubernetes

SPDK

- Host-managed array FTLs

- User-level device driver configurable/adapting to various workload

- Scalable design to manage array of NAND modules

- Leverage multi-core/multi-processor CPU to maximize parallelism

- Data center friendly orchestration utilizing Linux Containers and Kubernetes Ready design

CIRCUITBLVD
Change your memories. Change your life.

# CBOS™ Software Architecture

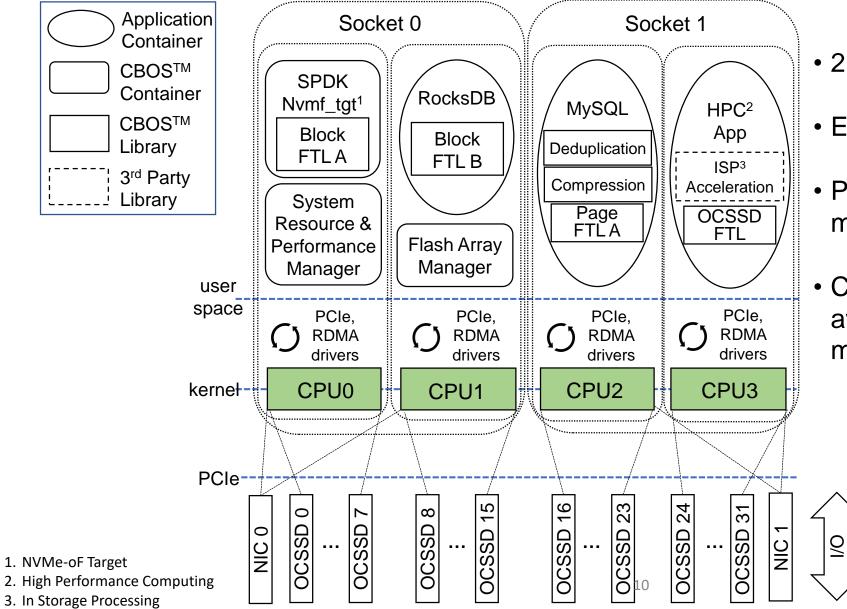## Container based storage and application software modules

- Host based Flash Array FTLs come with libraries that apps can pick and choose

\* OCSSD: Open-Channel SSD

CIRCUITBLVD
Change your memories. Change your life.

# CBOS™ Container Example

**Legend:**

- Application Container (oval)
- CBOS™ Container (rounded rectangle)
- CBOS™ Library (rectangle)
- 3rd Party Library (dashed rectangle)

**Socket 0**

SPDK Nvmf_tgt[1]
- Block FTL A
- System Resource & Performance Manager

RocksDB
- Block FTL B
- Flash Array Manager

**Socket 1**

MySQL
- Deduplication
- Compression
- Page FTL A

HPC[2] App
- ISP[3] Acceleration
- OCSSD FTL

**user space**

PCIe, RDMA drivers | PCIe, RDMA drivers | PCIe, RDMA drivers | PCIe, RDMA drivers

**kernel** — CPU0 | CPU1 | CPU2 | CPU3

**PCIe**

NIC 0 | OCSSD 0 ... OCSSD 7 | OCSSD 8 ... OCSSD 15 | OCSSD 16 ... OCSSD 23 | OCSSD 24 ... OCSSD 31 | NIC 1 | I/O

- 2 Socket, 4 Core System

- Each CPU core handles 8 OCSSD

- Per different application needs, matching FTL container is deployed

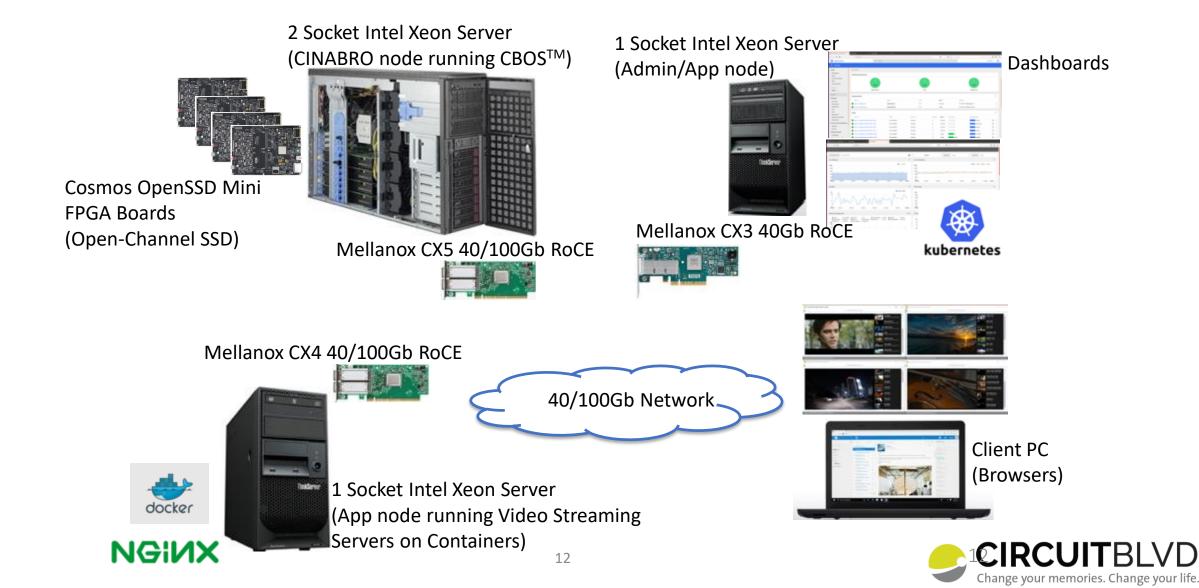- CBOS storage containers run on available CPU cores executing managerial tasks

1. NVMe-oF Target
2. High Performance Computing
3. In Storage Processing

10

CIRCUITBLVD
Change your memories. Change your life.

# Development Milestones

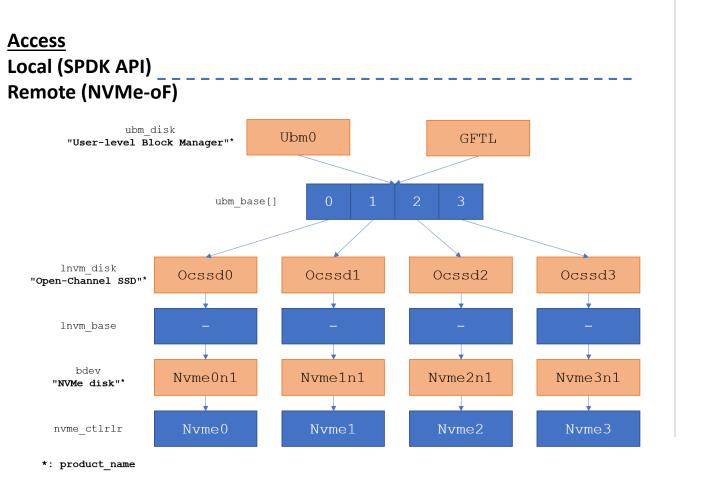| | Phase I (-Dec. '17) | Phase II (-Dec. '18) | Phase III (-Dec. '19) |
|---|---|---|---|
| **CINABRO™ Appliance** | **First Working Prototype**<br>• Commodity x86 server<br>• x4 Cosmos OpenSSD Mini PCIe cards<br>• Mellanox RDMA network cards | **Alpha**<br>• Commodity x86 server<br>• x8 Cosmos OpenSSD Ultra PCIe cards<br>• NVMe-oF network cards: TBD | **Beta**<br>• Customized PCIe fabrics<br>• Commodity CPU, DRAM, PCIe<br>• BMC: ready; Backup power: TBD |
| **CBBridge™** | **FPGA1**<br>• FPGA code with OCSSD compliant FW<br>• 16nm 2D MLC NAND w/ BCH | **FPGA2**<br>• RTL release: CBBridge™ SoC ready. Runs on FPGA<br>• 64L 3D TLC NAND w/ LDPC | **ASIC**<br>• SoC tape-out (mid '19)<br>• 96L 3D TLC/QLC NAND w/ LDPC |
| **CBOS™** | **Prototype release**<br>• NVMe-oF drive interface<br>• Baseline data path working:<br>• OCSSD pblk / lightnvm, NVMe-oF, SPDK/DPDK | **Alpha release**<br>• Host-based flash array FTLs<br>• Storage management layer<br>• Application plugins<br>• System resource & performance manager design complete | **Beta release**<br>• Core feature complete<br>• OpenStack compliant<br>• Data management beta<br>• System resource & performance Manager beta<br>• Out-of-Band management beta |
| **Open Source** | **SPDK contribution**<br>• Functions to help writing OCSSD access from user level<br>• Included in SPDK v17.10, v18.01 | **R&D version alpha**<br>• OpenSSD FPGA RTL codes v1.2<br>• Developer edition CBOS™ alpha: includes device drivers, user level libraries, and pilot apps | **R&D version beta**<br>• OpenSSD FPGA RTL codes v1.3<br>• Developer edition CBOS™ beta |

CIRCUITBLVD
Change your memories. Change your life.

# Current Prototype

## Multiple FPGA based OCSSDs running in our lab

2 Socket Intel Xeon Server
(CINABRO node running CBOS™)

1 Socket Intel Xeon Server
(Admin/App node)

Dashboards

Cosmos OpenSSD Mini
FPGA Boards
(Open-Channel SSD)

Mellanox CX5 40/100Gb RoCE

Mellanox CX3 40Gb RoCE

kubernetes

Mellanox CX4 40/100Gb RoCE

40/100Gb Network

Client PC
(Browsers)

docker

1 Socket Intel Xeon Server
(App node running Video Streaming
Servers on Containers)

NGINX

12

CIRCUITBLVD
Change your memories. Change your life.

# Prototype FTL Evaluations

Our 1st host FTLs, UBM and GFTL*, have been implemented in SPDK.

**Access**
**Local (SPDK API)**
**Remote (NVMe-oF)**



(1) Array with up to 4 units of FPGA OCSSDs

: showing reasonable performance for FPGA based SSDs

(2) Array with up to 24 units of OCSSD Qemu-nvme emulators

: used for qualitative test of multiple applications
: 24 copies of GFTL working correctly over same number of emulated OCSSDs

*: product_name

*GFTL is based on Hanyang University's Greedy FTL.*

13

# Prototype Applications

## End-to-end integration tests with multiple applications

| Containerized Video Server | Containerized RocksDB | SK telecom's AF Ceph |
|---|---|---|
| • OCSSD arrays host movie files and are exposed via SPDK nvmf_tgt containers<br>• Video servers run inside containers, made of Nginx web server with RTMP module | • Containerized RocksDB, both local and remote (over NVMe-oF)<br>• Local: SPDK's RocksDB plugin over CBOS™ GFTL<br>• Remote: SPDK's nvmf-tgt over CBOS™ GFTL | • All Flash Ceph is flash optimized version of Ceph.<br>• Initial data verification test over 4 FPGA OCSSDs ran successfully. |
| [8 concurrent video server example] | [24 RocksDB plugin per CPU core example] | [AF Ceph over 4 OCSSDs example] |

CIRCUITBLVD
Change your memories. Change your life.

# Summary

## Solution Benefits
- Flexibility to accommodate NAND generations from various vendors
- Adaptable to various Cloud Data Center network infrastructure
- Customizable SW architecture to meet ever-evolving cloud data center requirements

## Communities
- OpenSSD, OCSSD, SPDK: Our work has been integrated
- OCP: open to collaboration about making our hardware design available to the community

## Resources
- OpenSSD FPGA SSD available at: http://openssd.io
- SPDK OCSSD contributed codes: https://github.com/spdk/spdk.git  (SPDK v17.10, v18.01)
- CBOS$^{TM}$ development edition codes: TBD

**CIRCUIT**BLVD
Change your memories. Change your life.