



Flash Memory Summit

Raising QLC Reliability in All-Flash Arrays

Jeff Yang

Principal Engineer

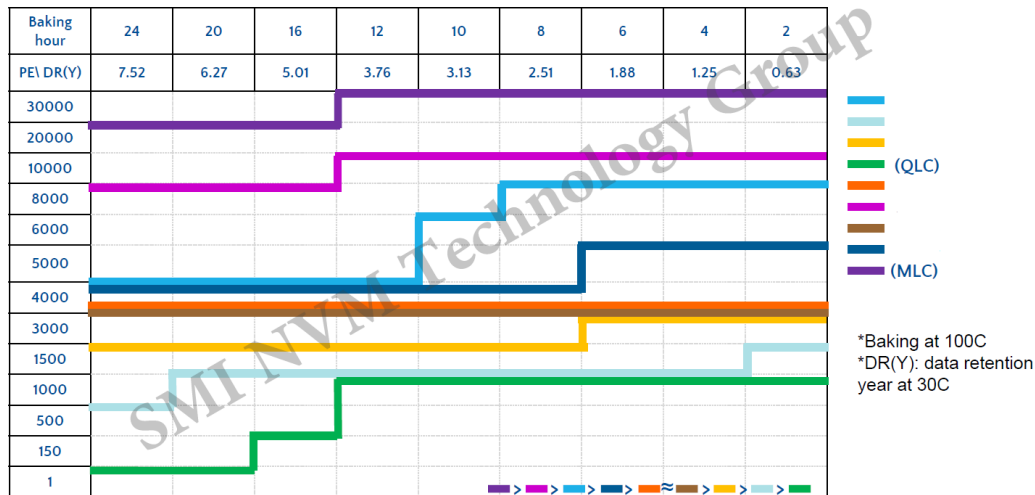
Storage Research Dept.

Silicon Motion, Inc.



QLC Characteristics

Endurance Rating (Estimation)

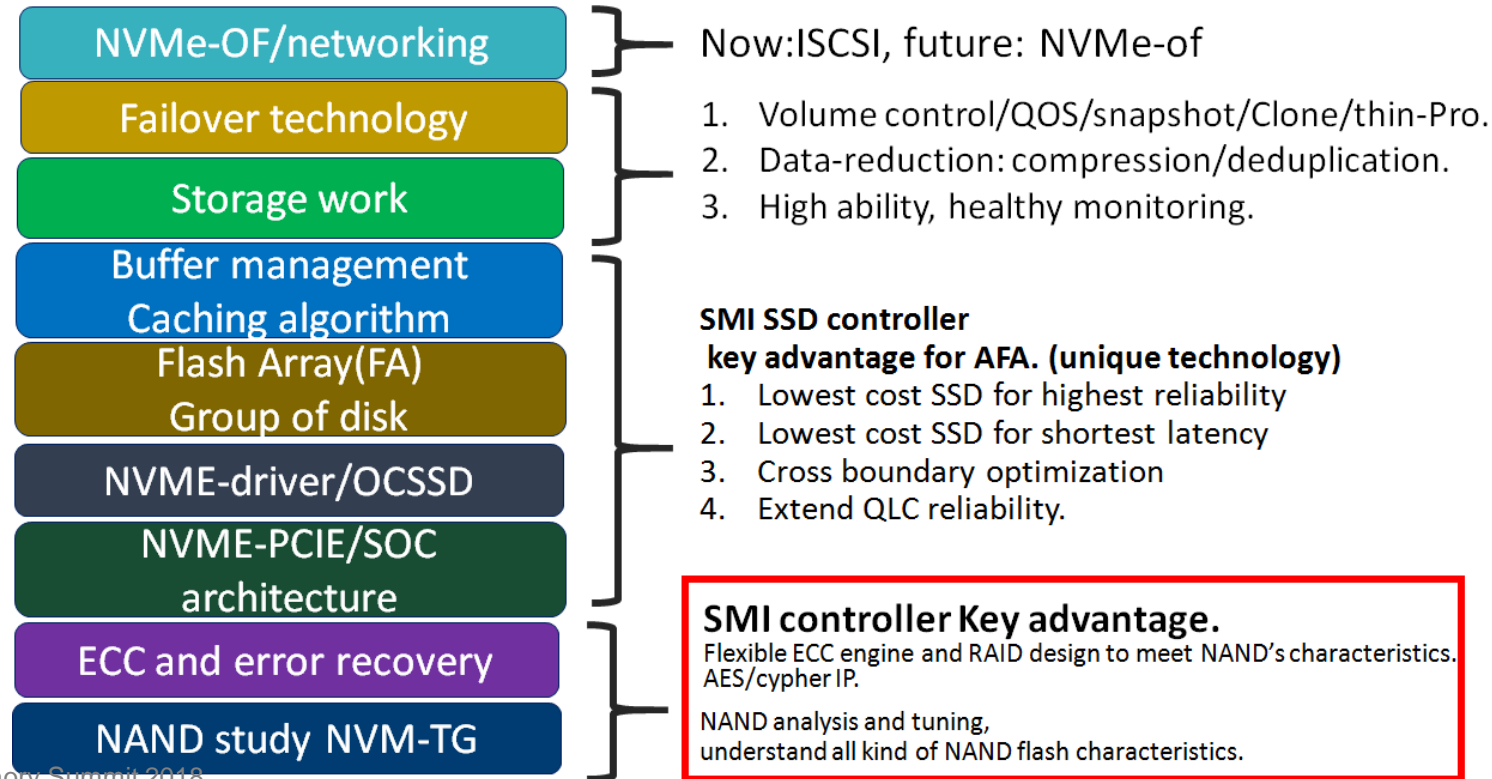


- QLC Endurance: 1~3K P/E.(limited DWDP)
- Enterprise TLC: 10K P/E.
- Average tPro per page is 3~4X to TLC.
- QLC Cost reduction ~25~30%.
- Ungraceful shutdown failure range on QLC is a conflict to the atomic write.
- SLC block from TLC or QLC flash are almost the same.

QLC is good enough for the SMB(Small and Medium-sized Business) AFA applications



Vertical integration for storage technology





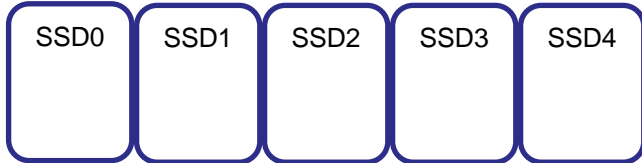
Flash Array

NVMe-OF/networking

Storage work

Buffer management
Caching algorithm

Flash Array (FA)
Group of disk



SLBA: System Logical block address

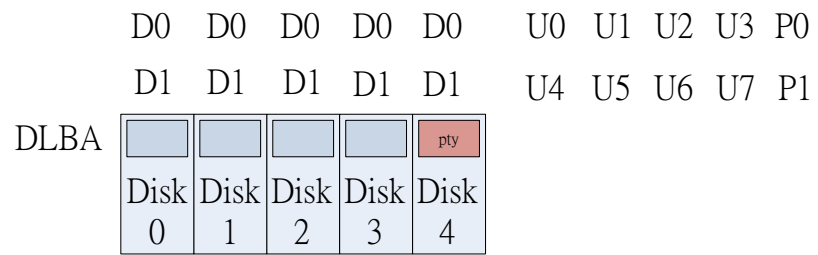
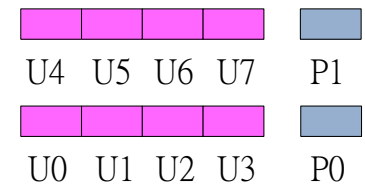
DLBA is SLBA's physical location

DLBA: Disk Logical block address

DLBA is Disk logical address

FPPA: Flash physical page address

RD U4, P1
 XOR P1 U4 → P1_tmp
 U4 XOR P1_tmp U4' → P1'
 WR U4', P1'





N SSDs with parity P, PQ, and PQR

- N + P: 1Write → 2Read + 2Write
→ Single parity: WAI = 2
- N + PQ: 1Write → 3 Read + 3Write
→ Double parity: WAI = 3
- N + PQR: 1Write → 4 Read + 4Write
→ Triple parity: WAI = 4
- The traditional method is not suitable for SSDs, because of the high WAI factor and consume the SSD's endurance faster.



Lower WAI RAID method on SSD

- Map the SLBA to the DLBA.
- Generate two parity on the same DLBA cross different SSDs to provide the protection.
- Flash array software layer maintain a lookup table.
- When writing the existed data into the Flash array:
 - ➔ Existed RAID link will not be changed.
 - ➔ Write the data to the new location with new RAID link.

L2P table

U1	D ₀₀
U2	D ₁₀
U3	D ₂₀
U4	D ₀₁
U5	D ₁₁
U6	D ₂₁
U7	D ₀₂
U8	D ₁₂
U9	D ₂₂

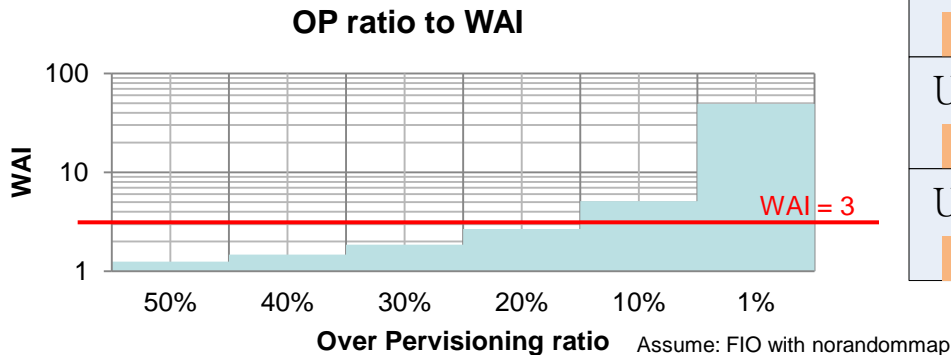
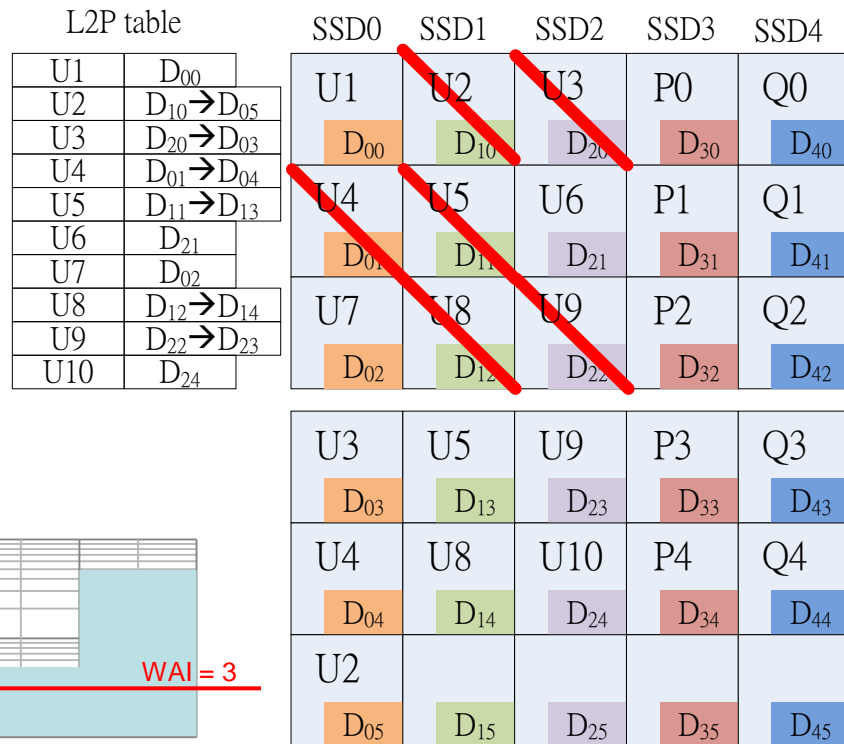
SSD0	SSD1	SSD2	SSD3	SSD4
U1 D ₀₀	U2 D ₁₀	U3 D ₂₀	P0 D ₃₀	Q0 D ₄₀
U4 D ₀₁	U5 D ₁₁	U6 D ₂₁	P1 D ₃₁	Q1 D ₄₁
U7 D ₀₂	U8 D ₁₂	U9 D ₂₂	P2 D ₃₂	Q2 D ₄₂

D ₀₃	D ₁₃	D ₂₃	D ₃₃	D ₄₃
D ₀₄	D ₁₄	D ₂₄	D ₃₄	D ₄₄
D ₀₅	D ₁₅	D ₂₅	D ₃₅	D ₄₅



Lower WAI RAID method on SSD

- Write the data to the new location with new RAID link
- Invalidate the old location and update the L2P table.
- The GC work will be applied.
- The WAI will be related to the Overprovisioning (OP) ratio.





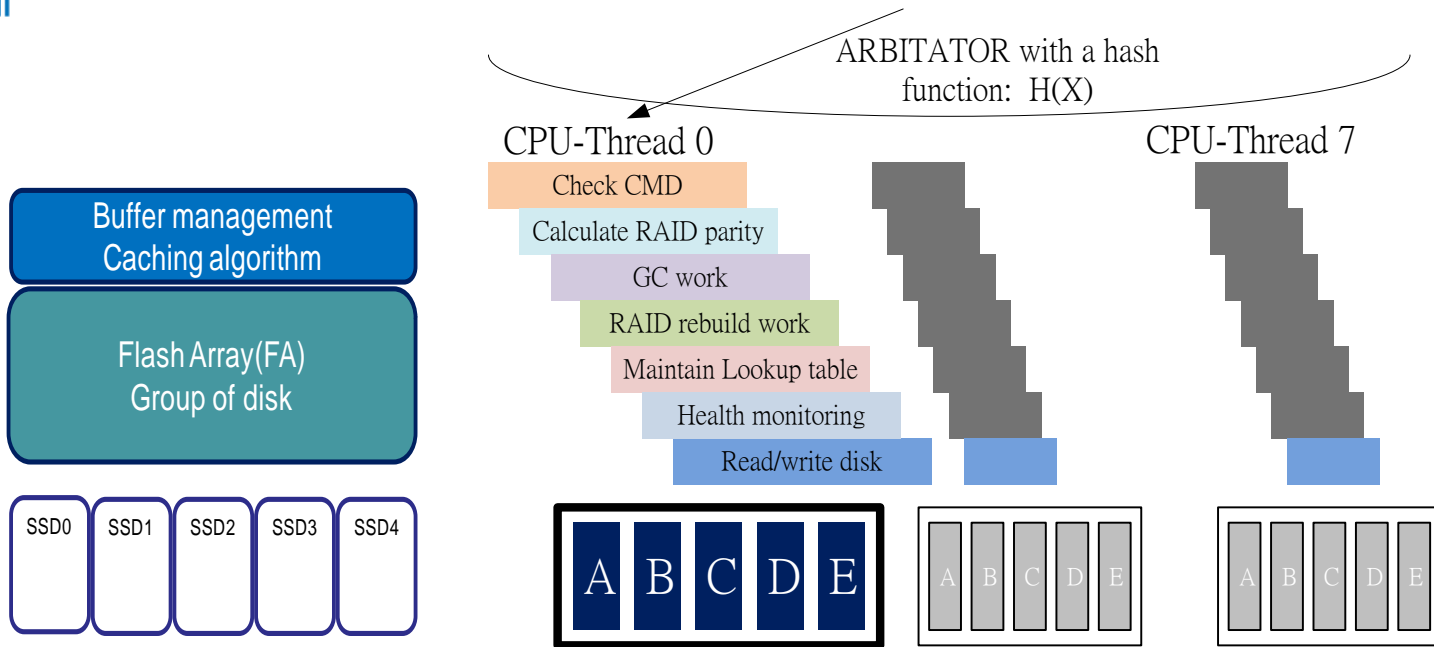
Single SSD Read/Write behavior

- The SSD awareness RAID will construct a remapping table in host side.
- For the single SSD point of view, the access behavior become the sequential write. (DLBA is a sequential write.)
- When issuing a read CMD on SLBA, host will redirect to DLBA by using the mapping table.

SSD0	SSD1	SSD2	SSD3	SSD4
D ₀₀	D ₁₀	D ₂₀	D ₃₀	D ₄₀
D ₀₁	D ₁₁	D ₂₁	D ₃₁	D ₄₁
D ₀₂	D ₁₂	D ₂₂	D ₃₂	D ₄₂
D ₀₃	D ₁₃	D ₂₃	D ₃₃	D ₄₃
D ₀₄	D ₁₄	D ₂₄	D ₃₄	D ₄₄
D ₀₅	D ₁₅	D ₂₅	D ₃₅	D ₄₅
D ₀₆	D ₁₆	D ₂₆	D ₃₆	D ₄₆
D ₀₇	D ₁₇	D ₂₇	D ₃₇	D ₄₇
D ₀₈	D ₁₈	D ₂₈	D ₃₈	D ₄₈



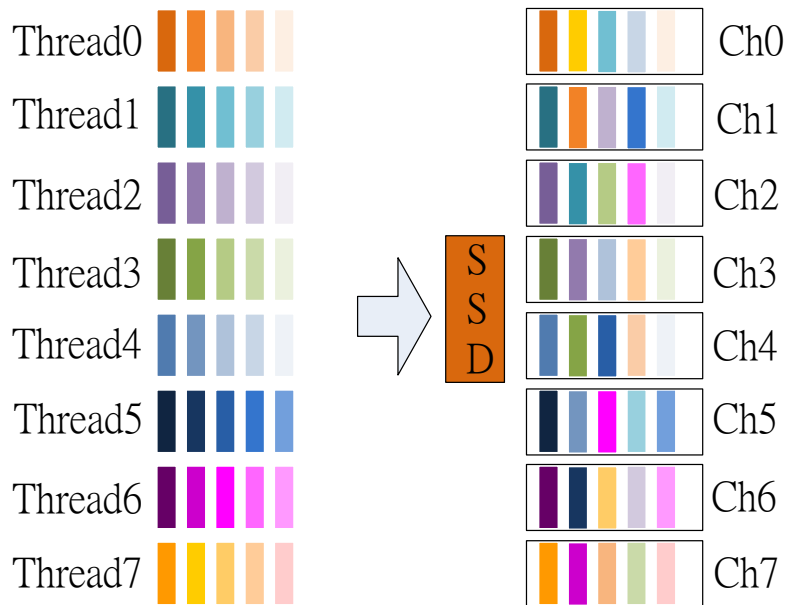
Multi-thread operation scenario



- Read/write disk using single thread will cause context switch problems.
- Single disk will need to handle the requests from different threads.



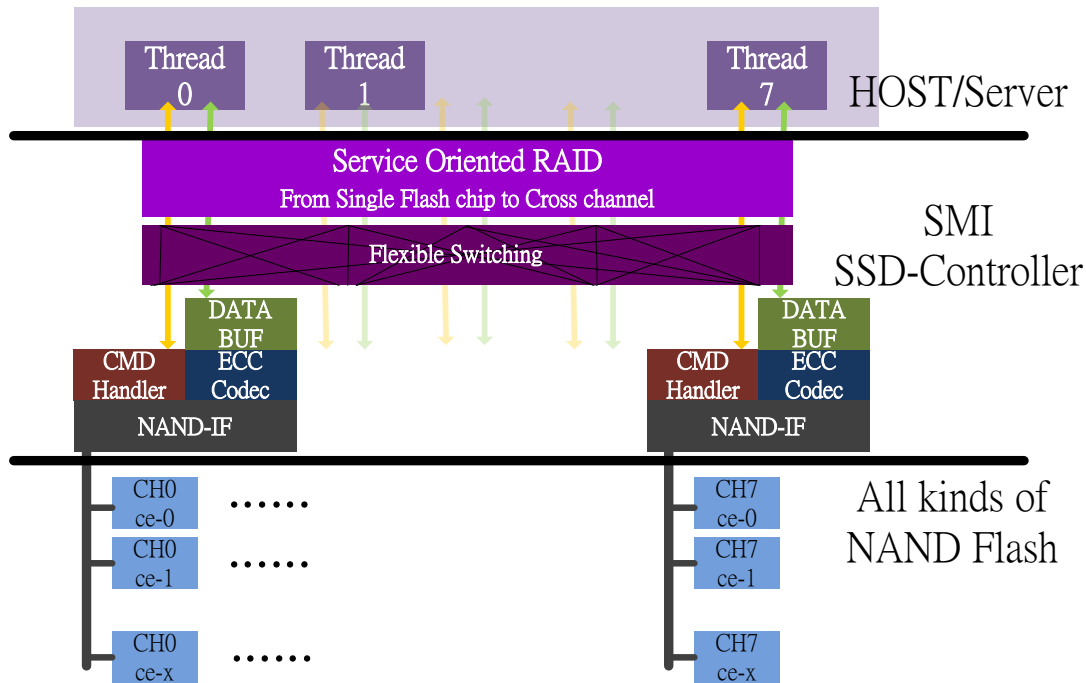
Access mixing from thread to SSD channel



- All the threads issue the access to the same SSD.
- It becomes the random access behavior.
- Huge **DRAM** for SSD device mapping table and another OP-ratio are required.
- Powerful SSD device CPU for GC work.
- Larger **Capacitance** for ungraceful shutdown handling.
- Enterprise SSD: OP=20%, WAI = ~3.
- Much more expensive enterprise SSD.



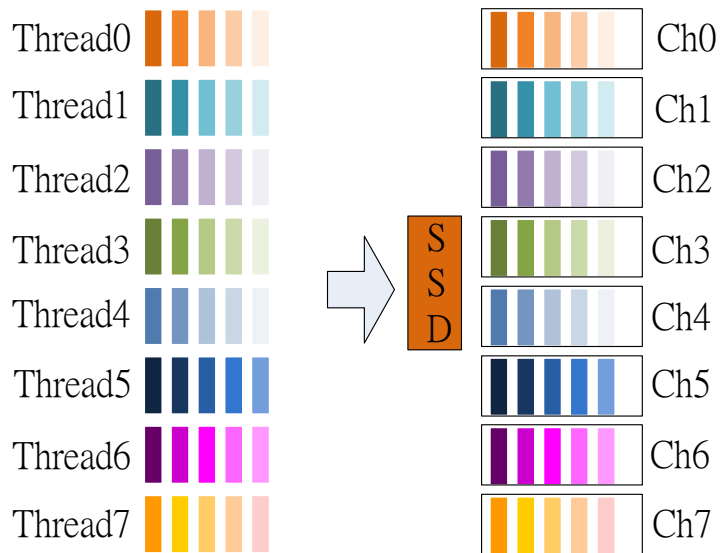
A flexible controller to solve problem. Service oriented SSD



- Flexible switching provides several different types of NAND groups for applications.
- The service oriented RAID is configurable for different types of NAND groups and different types of NAND failure behavior.
- Each Channel becomes sequential program and erase.
- Use SLC as caching buffer on ungraceful shutdown handling.



Every thread gets its own NAND Flash.



- Dedicated NANDs for dedicated threads respectively.
- The simple mapping **removes DRAM** requirement in SSD device.
- Remove GC work from device SSD
- Data will write into SLC first, **remove Capacitance** in ungraceful shutdown handling flow.
- Cost efficient SSD.
- **Sequential Write/Erase** behavior is the perfect match for QLC.



Comparison on N + PQR (triple RAID-parity)

	Traditional Flash array RAID	SSD awareness RAID	Vertical integration RAID
Host Flash array (SLBA to DLBA)	1W = 4R + 4W. OP = ~0% WAI = 4	OP = 20% WAI = 3	OP = 20% WAI = 3
SSD device (DLBA to FPPA)	Enterprise SSD. OP = 20%. WAI = 3	Enterprise SSD. OP = 20%. WAI = 3	Service oriented SSD. OP = ~0%. WAI = ~1
Overall (SLBA to FPPA)	WAI = 12 OP = 20% (additional read latency)	WAI = 9 OP = ~36%	WAI = 3. OP = 20%

- Reduce the overall WAI and Over Provisioning will increase the life time of QLC



Conclusion

- SSD controller is a key to connect the NAND to applications.
- Both reliability and efficiency will be improved by controller.

Hard decoding

Moving Read (Using read-retry table)

Soft-decoding (Iterative decoding)

HRE aware iterative decoding

RAID protection (addition parity)

Cross disk RAID protection

Remote backup