



Flash Memory Summit

ASIC/Merchant Chip-Based Flash Controllers

Dr. Jeff Yang

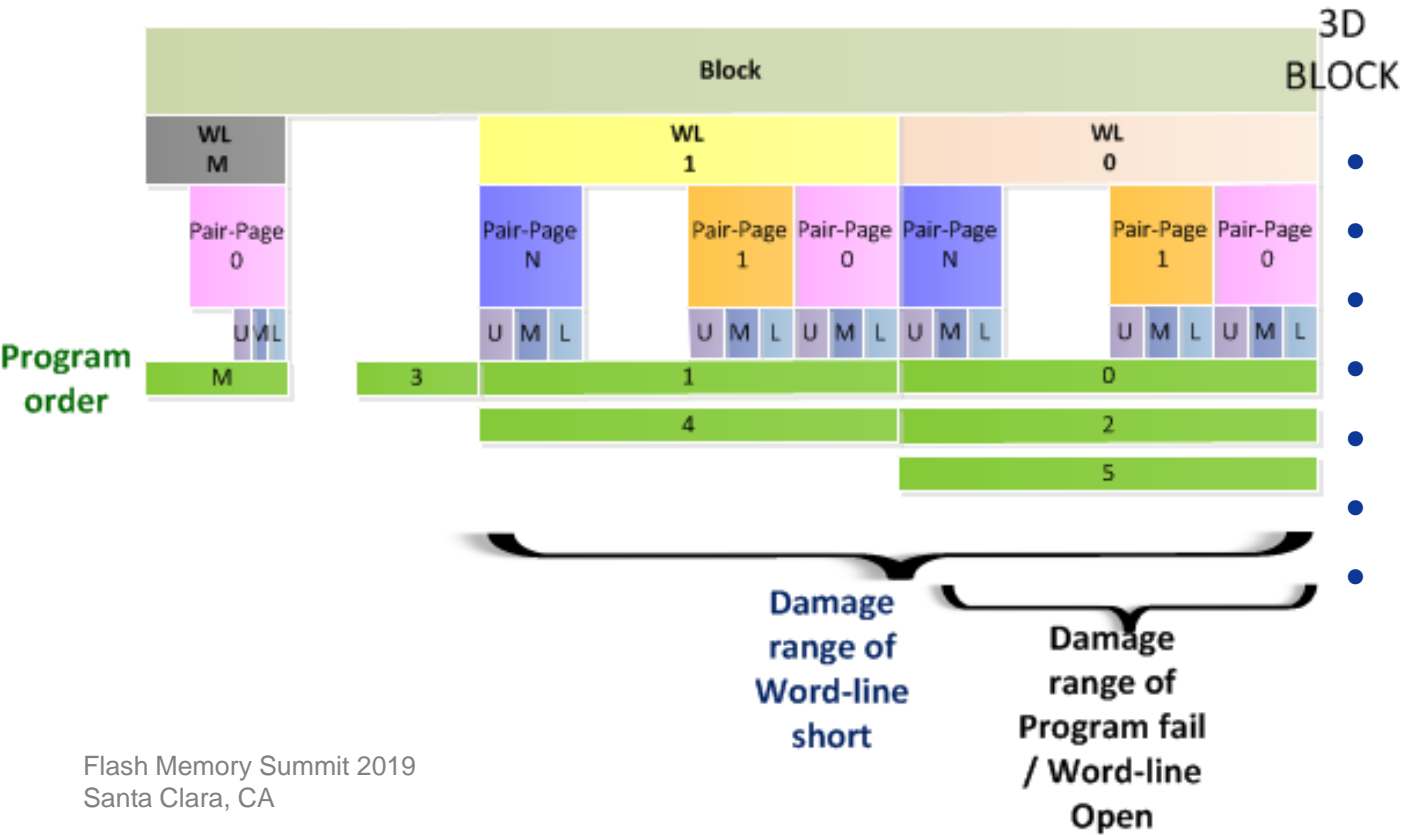
Principal Engineer

Storage Research Dept.

Silicon Motion, Inc.



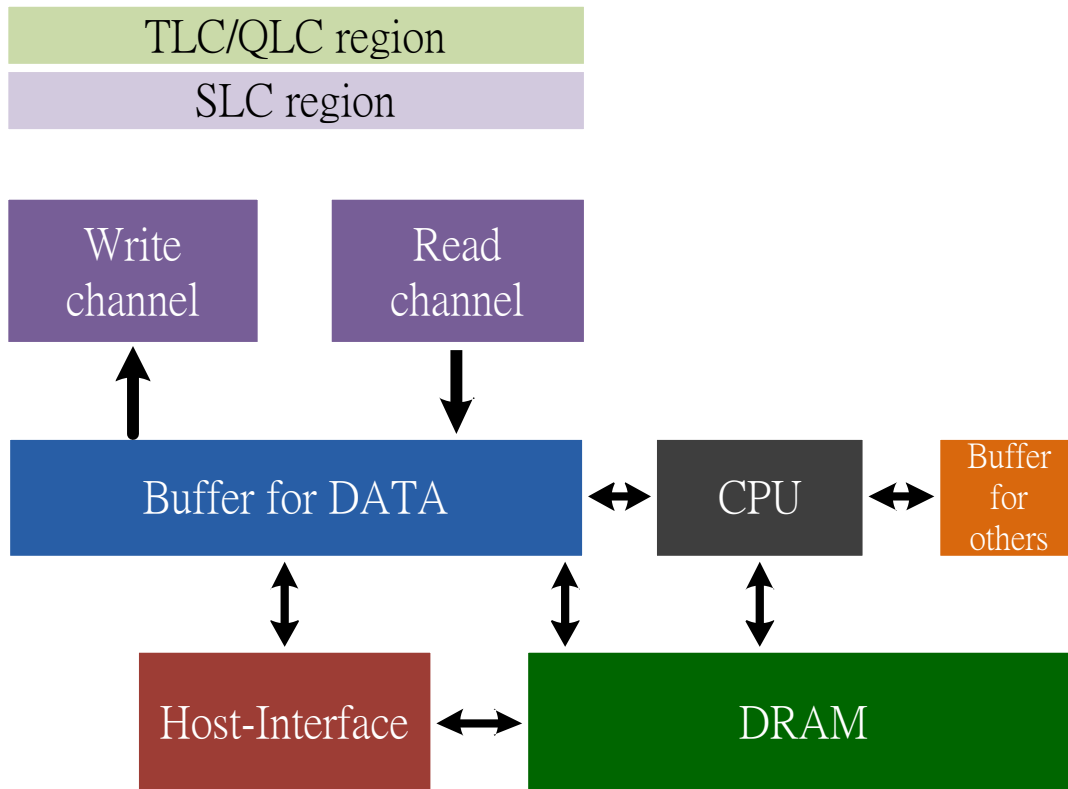
NAND usage/failure/defect



- Program order.
- Program failure.
- Ungraceful shutdown.
- Voltage detection.
- Uncompleted WL
- Uncompleted block
- SLC/TLC/QLC usage

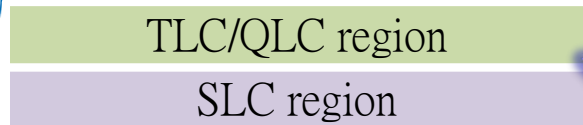


Basic Architecture of SSD controller

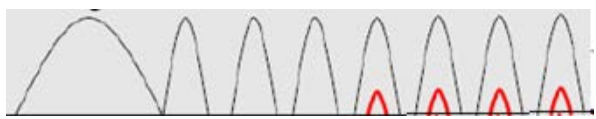
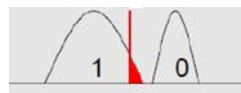
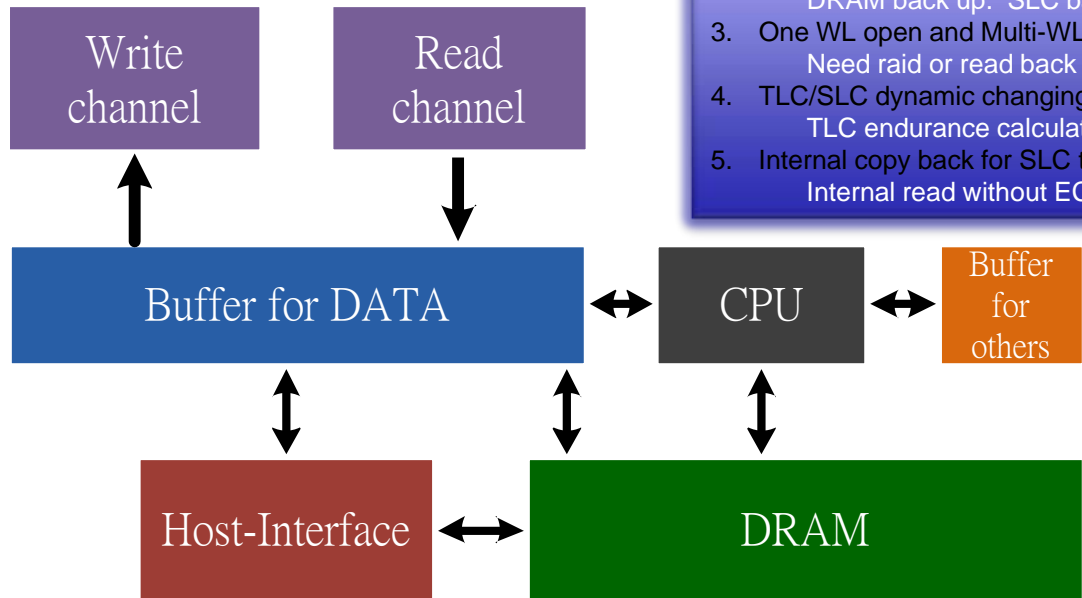




Different TLC/QLC reliability issue

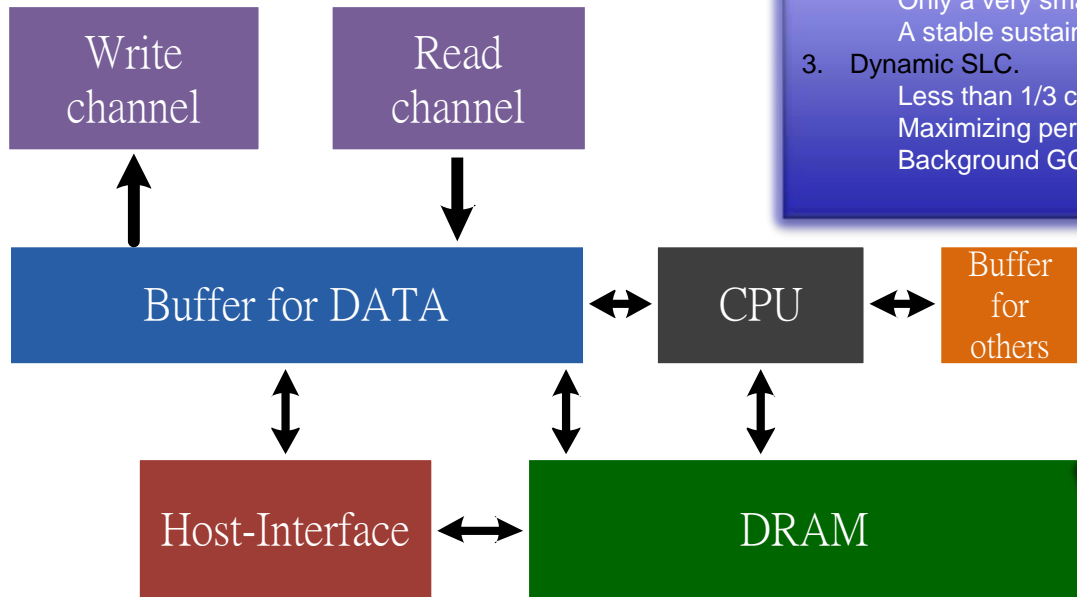
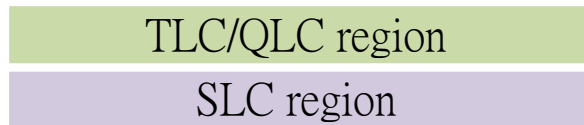


- 3D. TLC/QLC
One-pass program, two-pass program, multi-pass program.
- Program failure protection flow.
Program failure range.
Read back data from flash cache buffer.
DRAM back up. SLC back up, SRAM backup.
- One WL open and Multi-WL short.
Need raid or read back check after program.
- TLC/SLC dynamic changing usage.
TLC endurance calculation issue.
- Internal copy back for SLC to TLC
Internal read without ECC correction





Application combinations on TLC/QLC



1. SLC caching.
Data always write into SLC. A Fixed portion of SLC regard as a cache buffer.
Performance boost on SLC caching.
Background GC to TLC block.
2. TLC/QLC direct.
Only a very small portion for system usage and small random write data.
A stable sustained write performance.
3. Dynamic SLC.
Less than 1/3 capacity threshold, using SLC.
Maximizing performance boosting period.
Background GC to TLC block.

1. Full size DRAM
For host data write caching. Full lookup table.
2. Non-DRAM
Extremely low cost.
Optimize for user experience.
More system info access from SLC blocks.
3. Small DRAM.
Full lookup table on external buffer
No host data buffer.



Challenge: Support all combinations and cost efficiency

3D TLC/QLC

One-pass

Two-pass

Multi-pass

SLC usage

SLC caching

TLC direct

SLC/(TLC/QLC)
dynamic

External buffer

Full DRAM

Non-DRAM

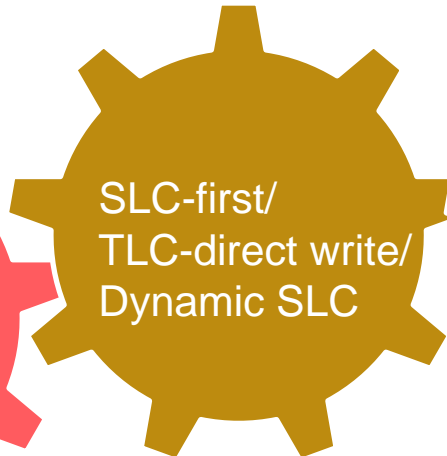
Partial DRAM



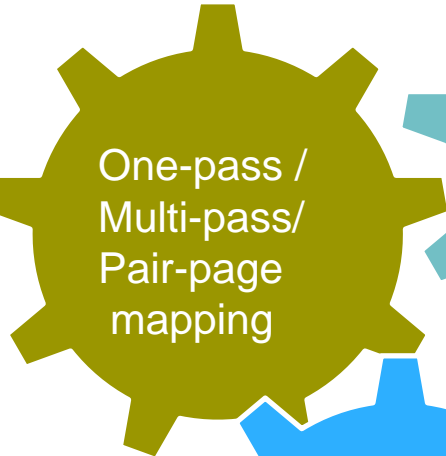
Flash Memory Summit



DRAM/
DRAM-less/
Small DRAM



SLC-first/
TLC-direct write/
Dynamic SLC

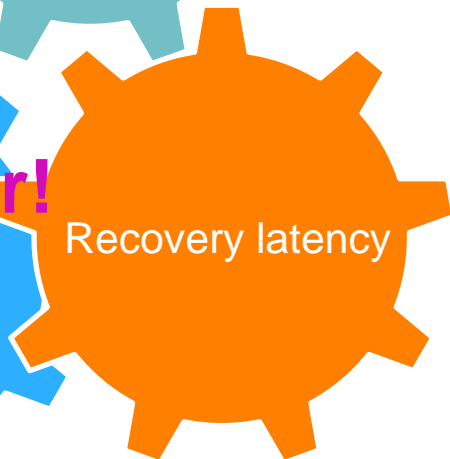


One-pass /
Multi-pass/
Pair-page
mapping



WL to WL short
Failure range

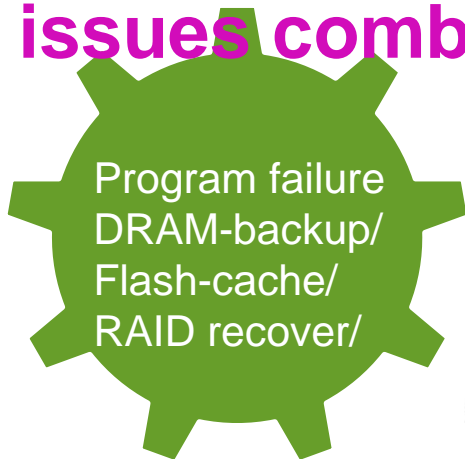
All the issues combine together!



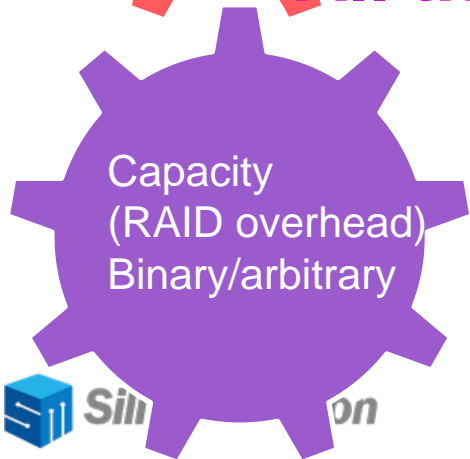
Recovery latency



WL open
Failure range



Program failure
DRAM-backup/
Flash-cache/
RAID recover/



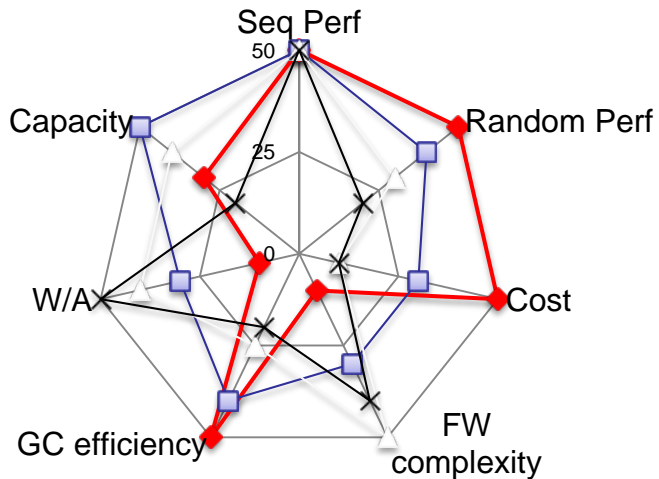
Capacity
(RAID overhead)
Binary/arbitrary





Product comparison

- ◆ Full DRAM
- Partial DRAM
- △ None DRAM (HMB)
- ✕ None DRAM

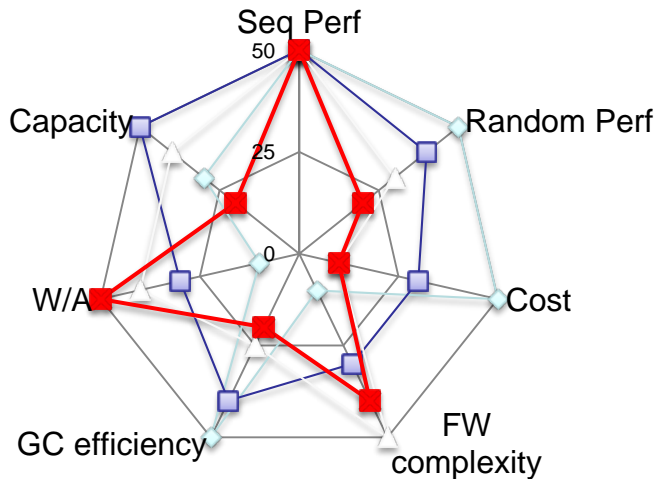


Full DRAM	Partial DRAM	None DRAM	None DRAM (HMB)
<p>Pros</p> <ul style="list-style-type: none"> • High Perf • GC efficiency • Low W/A • Low complexity 	<p>Pros</p> <ul style="list-style-type: none"> • GC efficiency • Capacity limited by L2P address bit (8TB) 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Perf still good for some user behavior 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Random perf
<p>Cons</p> <ul style="list-style-type: none"> • Capacity limitation (2TB) • DRAM bandwidth limitation • Power consumption • Cost 	<p>Cons</p> <ul style="list-style-type: none"> • Med Random Perf (DRAM L2P buffer) • Med W/A 	<p>Cons</p> <ul style="list-style-type: none"> • High W/A • Low Random Perf (SRAM L2P buffer) • Program Failure • High capacity support 	<p>Cons</p> <ul style="list-style-type: none"> • Same as "None DRAM" • More complex than "None DRAM"



Product comparison

- ◆ Full DRAM
- Partial DRAM
- △ None DRAM (HMB)
- None DRAM

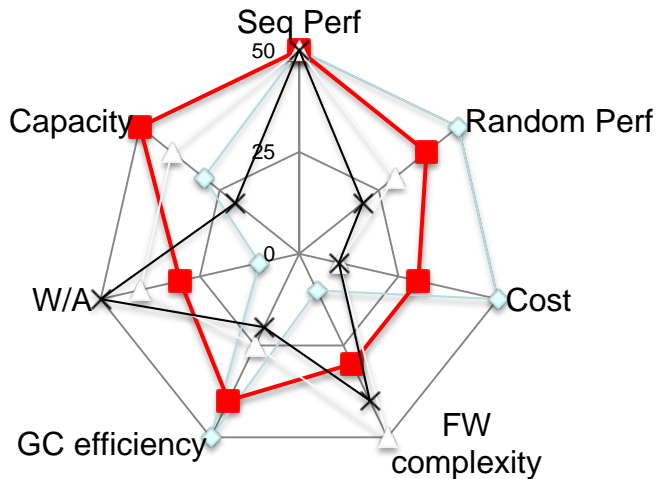


Full DRAM	Partial DRAM	None DRAM	None DRAM (HMB)
<p>Pros</p> <ul style="list-style-type: none"> • High Perf • GC efficiency • Low W/A • Low complexity 	<p>Pros</p> <ul style="list-style-type: none"> • GC efficiency • Capacity limited by L2P address bit (8TB) 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Perf still good for some user behavior 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Random perf
<p>Cons</p> <ul style="list-style-type: none"> • Capacity limitation (2TB) • DRAM bandwidth limitation • Power consumption • Cost 	<p>Cons</p> <ul style="list-style-type: none"> • Med Random Perf (DRAM L2P buffer) • Med W/A 	<p>Cons</p> <ul style="list-style-type: none"> • High W/A • Low Random Perf (SRAM L2P buffer) • Program Failure • High capacity support 	<p>Cons</p> <ul style="list-style-type: none"> • Same as "None DRAM" • More complex than "None DRAM"



Product comparison

- ◆ Full DRAM
- Partial DRAM
- △ None DRAM (HMB)
- ✕ None DRAM

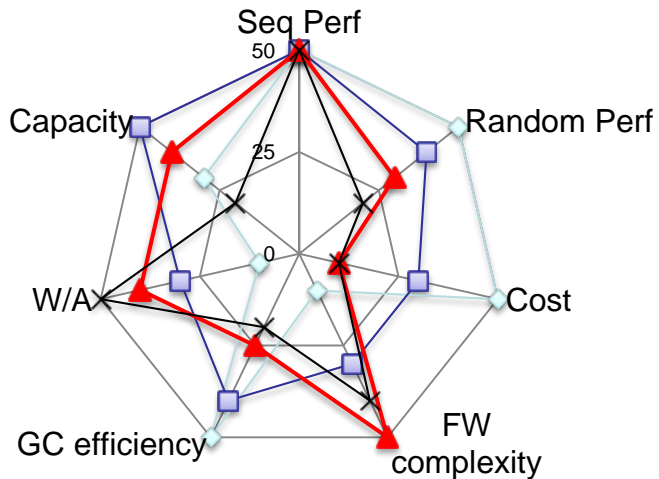


Full DRAM	Partial DRAM	None DRAM	None DRAM (HMB)
<p>Pros</p> <ul style="list-style-type: none"> • High Perf • GC efficiency • Low W/A • Low complexity 	<p>Pros</p> <ul style="list-style-type: none"> • GC efficiency • Capacity limited by L2P address bit (8TB) 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Perf still good for some user behavior 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Random perf
<p>Cons</p> <ul style="list-style-type: none"> • Capacity limitation (2TB) • DRAM bandwidth limitation • Power consumption • Cost 	<p>Cons</p> <ul style="list-style-type: none"> • Med Random Perf (DRAM L2P buffer) • Med W/A 	<p>Cons</p> <ul style="list-style-type: none"> • High W/A • Low Random Perf (SRAM L2P buffer) • Program Failure • High capacity support 	<p>Cons</p> <ul style="list-style-type: none"> • Same as "None DRAM" • More complex than "None DRAM"



Product comparison

- ◆ Full DRAM
- Partial DRAM
- ▲ None DRAM (HMB)
- ✕ None DRAM



Full DRAM	Partial DRAM	None DRAM	None DRAM (HMB)
<p>Pros</p> <ul style="list-style-type: none"> • High Perf • GC efficiency • Low W/A • Low complexity 	<p>Pros</p> <ul style="list-style-type: none"> • GC efficiency • Capacity limited by L2P address bit (8TB) 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Perf still good for some user behavior 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Random perf
<p>Cons</p> <ul style="list-style-type: none"> • Capacity limitation (2TB) • DRAM bandwidth limitation • Power consumption • Cost 	<p>Cons</p> <ul style="list-style-type: none"> • Med Random Perf (DRAM L2P buffer) • Med W/A 	<p>Cons</p> <ul style="list-style-type: none"> • High W/A • Low Random Perf (SRAM L2P buffer) • Program Failure • High capacity support 	<p>Cons</p> <ul style="list-style-type: none"> • Same as "None DRAM" • More complex than "None DRAM"



MLC

- Pros:
- Aligned to OS LBA size
 - Ease to do GC
 - Good perf
- Cons:
- 1. High Cost

4K pa
mapp

- Pros:
- Good Perf (SLC) for FOB
 - Med sustained perf
- Cons:
- High latency during low power transition
 - Lot data loss in SPOR

TLC

- Pros:
- Good Perf (SLC) for FOB
 - Reliability
- Cons:
- 1. Low sustained perf

SLC first

SLC
Caching

Dynamic
SLC

SLC-To-TLC

Internal
CopyBack

External
GC

Idle
eviction

External
GC

- Pros:
- Good perf during GC
- Cons:
- Err bits accumulated
 - High W/A

- Pros:
- Low W/A
- Cons:
- Low perf during GC

- Pros:
- Better user experience
- Cons:
- Low sustained perf
 - High GC effort

3D TLC

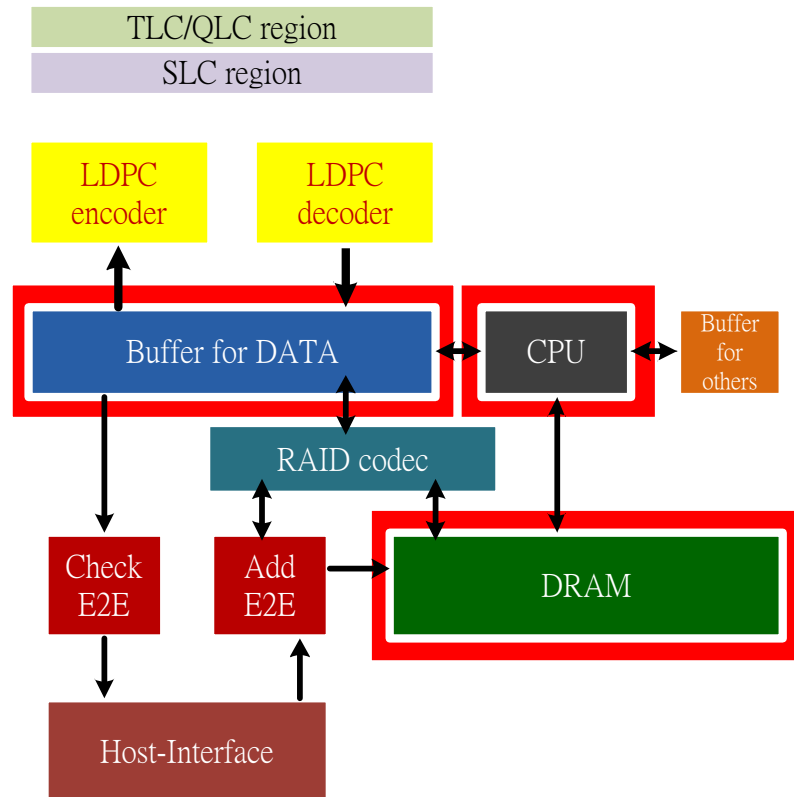
TLC direct

- Pros:
- Good sustained perf
- Cons:
- Data loss in SPOR



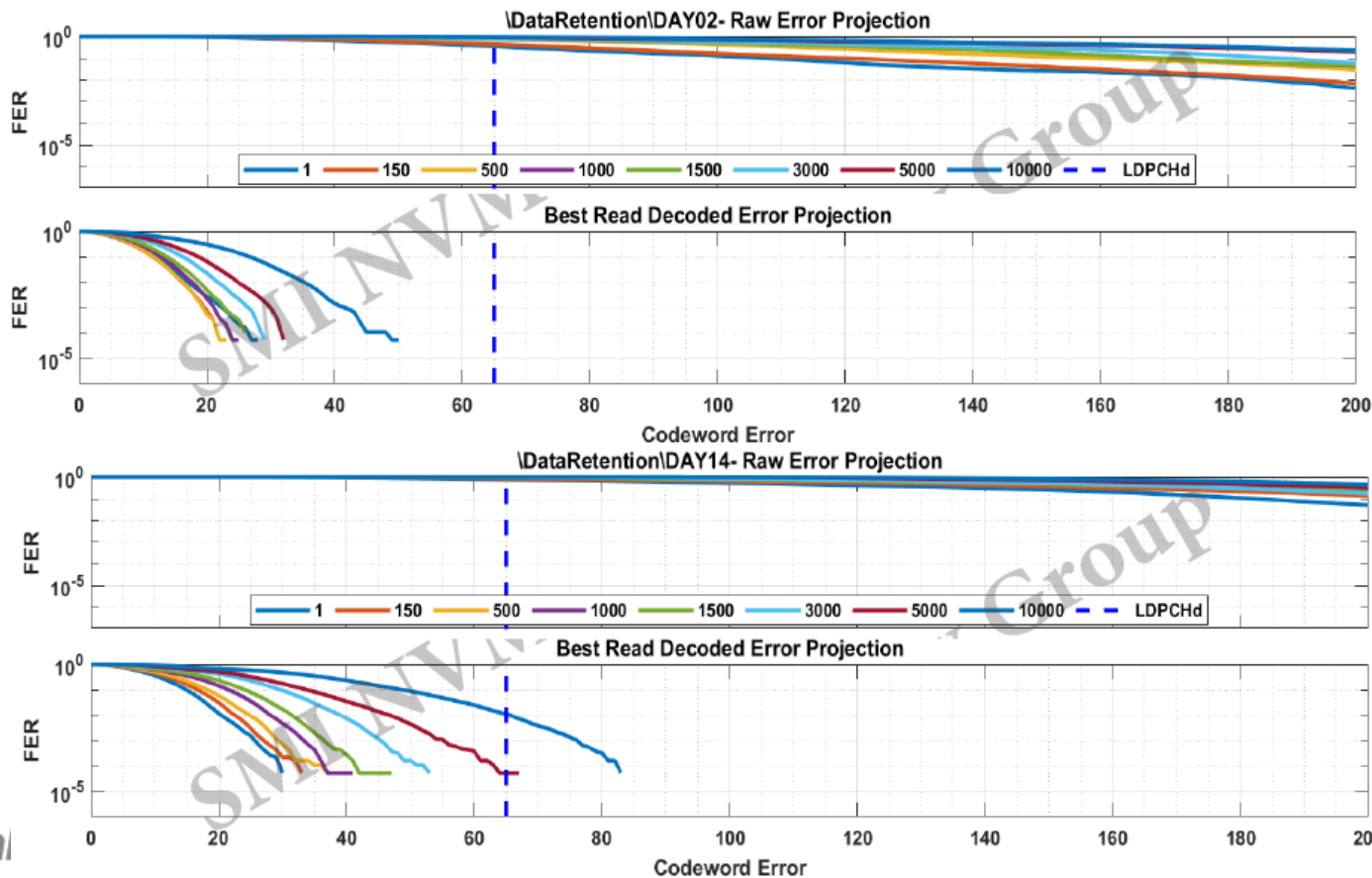
Enhanced reliability on SSD controller

- ECC to NAND E2E
- Stronger LDPC engine. Increase to 4KB LDPC.
- SRAM-ECC.
- DRAM-ECC.
- Host interface E2E



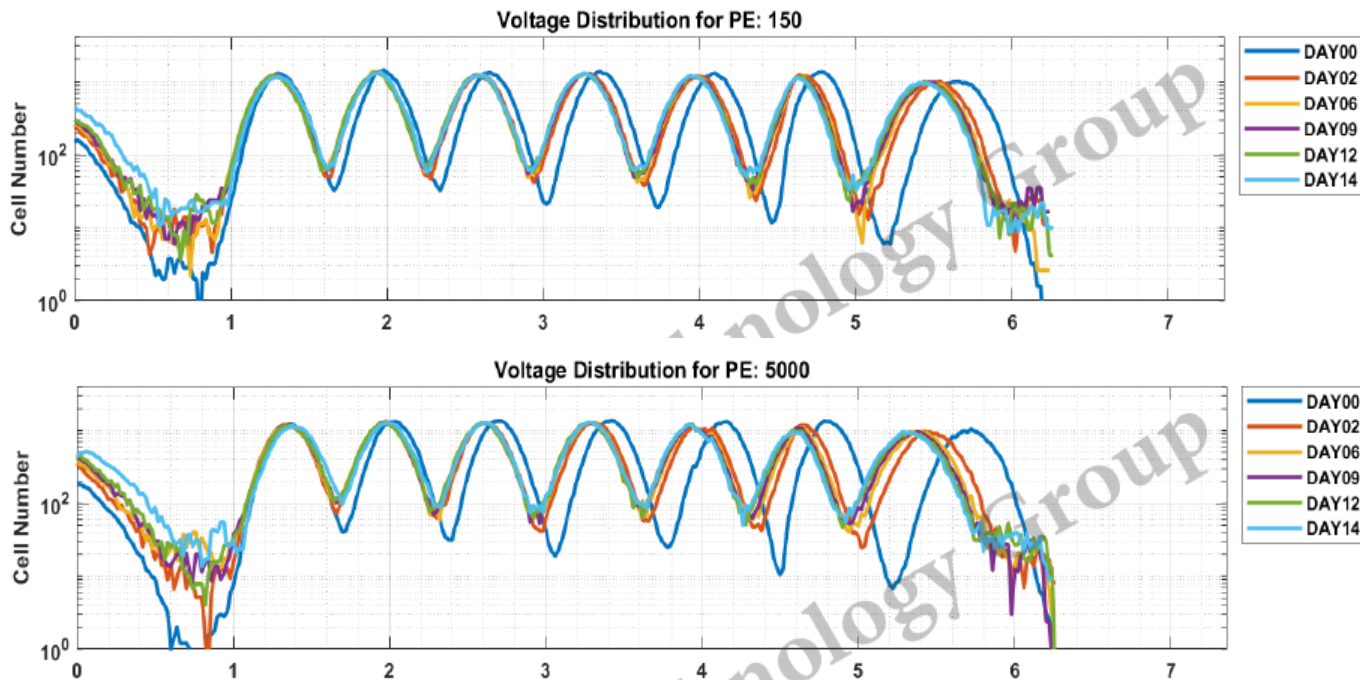


Low temp data retention



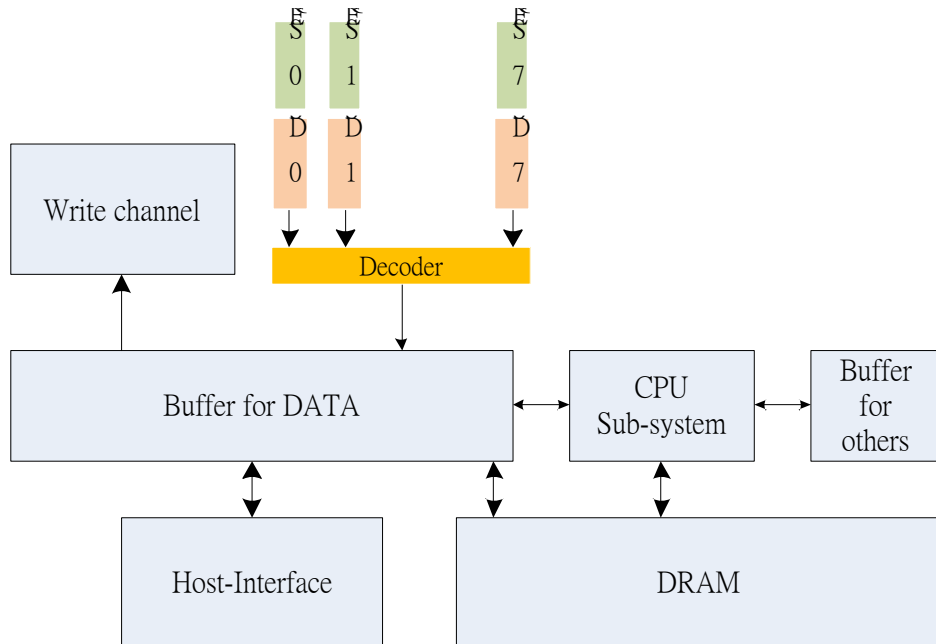
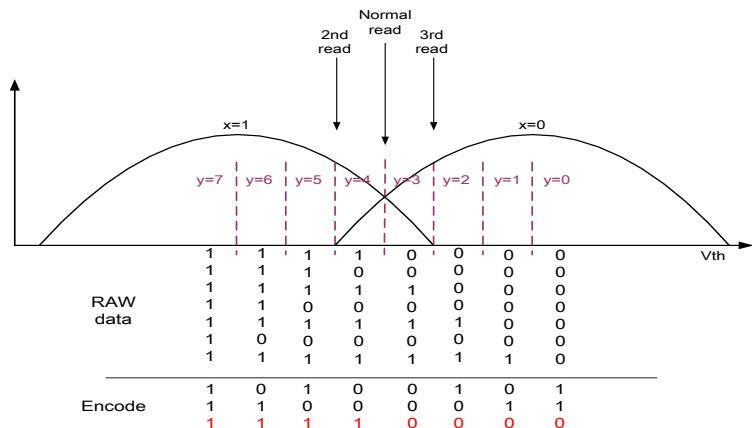


Low temp data retention issue.





Soft-information interface



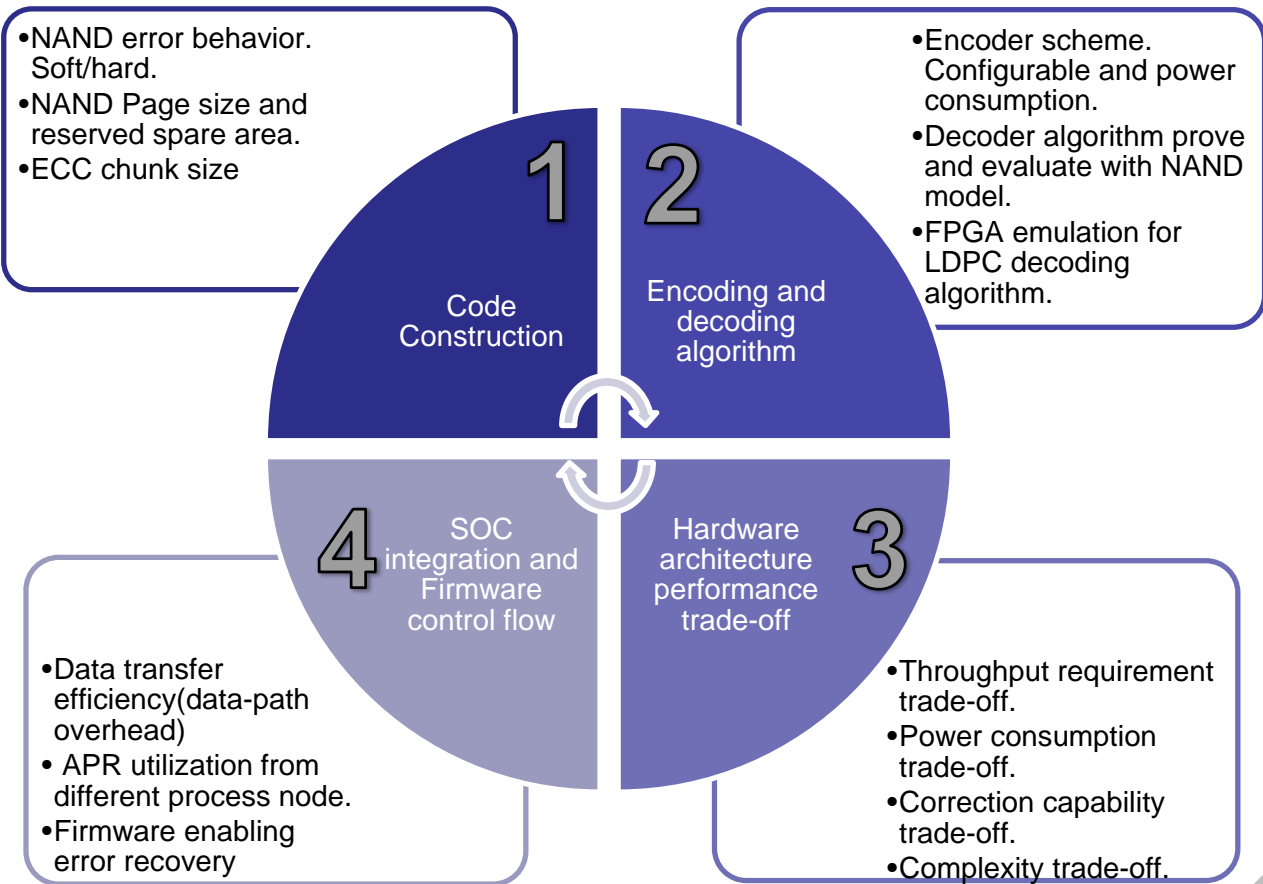
- In order to provide better decoder's correction capability, using the soft-info to get more reliability bits.
- NAND interface support .
 - Traditional read/retry interface.
 - Direct soft-info interface.



ECC design loop related to NAND characteristics.

Flash Memory Summit

- Keep improving the LDPC performance.
- For higher throughput 8~16GB/sec, we may go back to step1.
- After 16nm process, the design iteration depth will from code-construction to trial APR.
- EX: Find the Routing congestion issue in step 4, it may need to solve from step1.
- 12nm process is another new story on the LDPC engine design



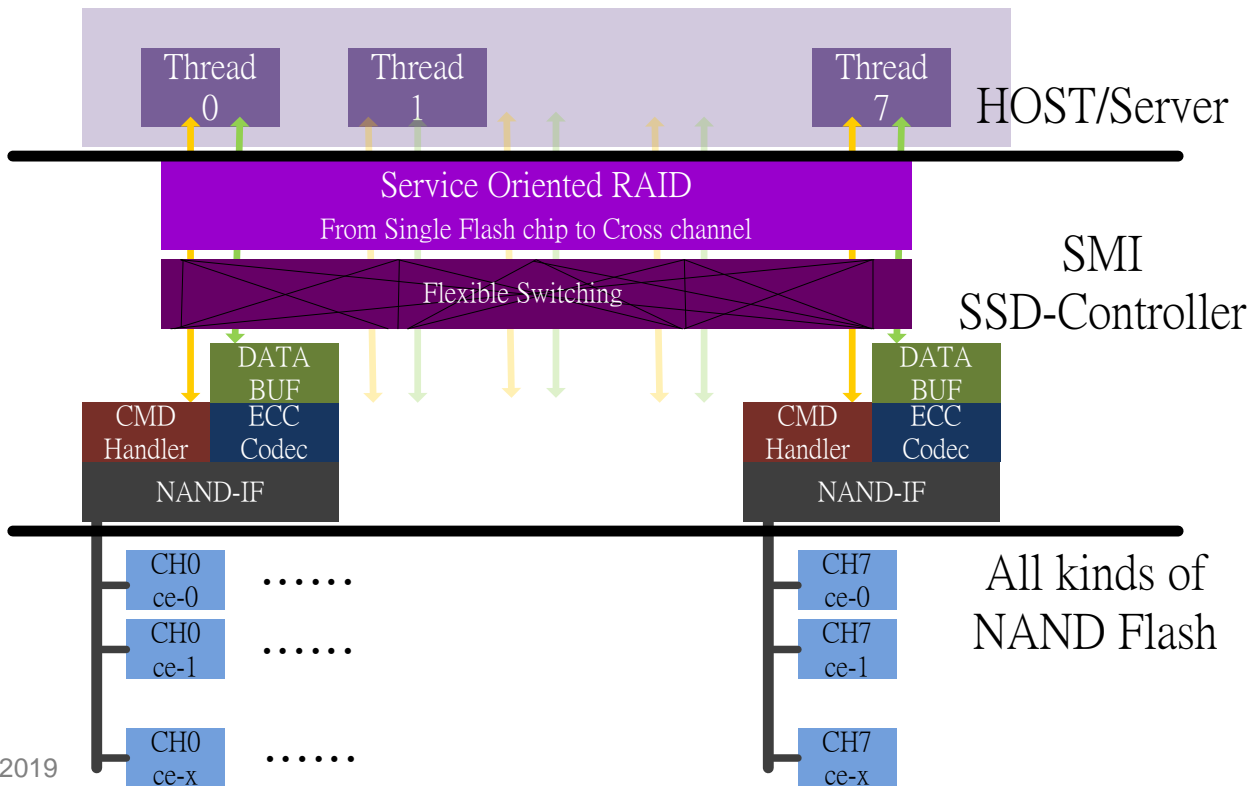


12nm process impact the and SOC design

- SRAM power consumption and the routing congestion.
- Single NAND channel will have higher than 1.2GB/sec throughput.
- 4-channel for Client SSD, 8-channel for enterprise and 16-channel for Data-center.
- Latency consistency. Decoding latency may blockage the chunk decoding between channels.
 - Program suspend for read.
 - Erase suspend for read.
- Controller buffer management for the read/write mixed behavior.

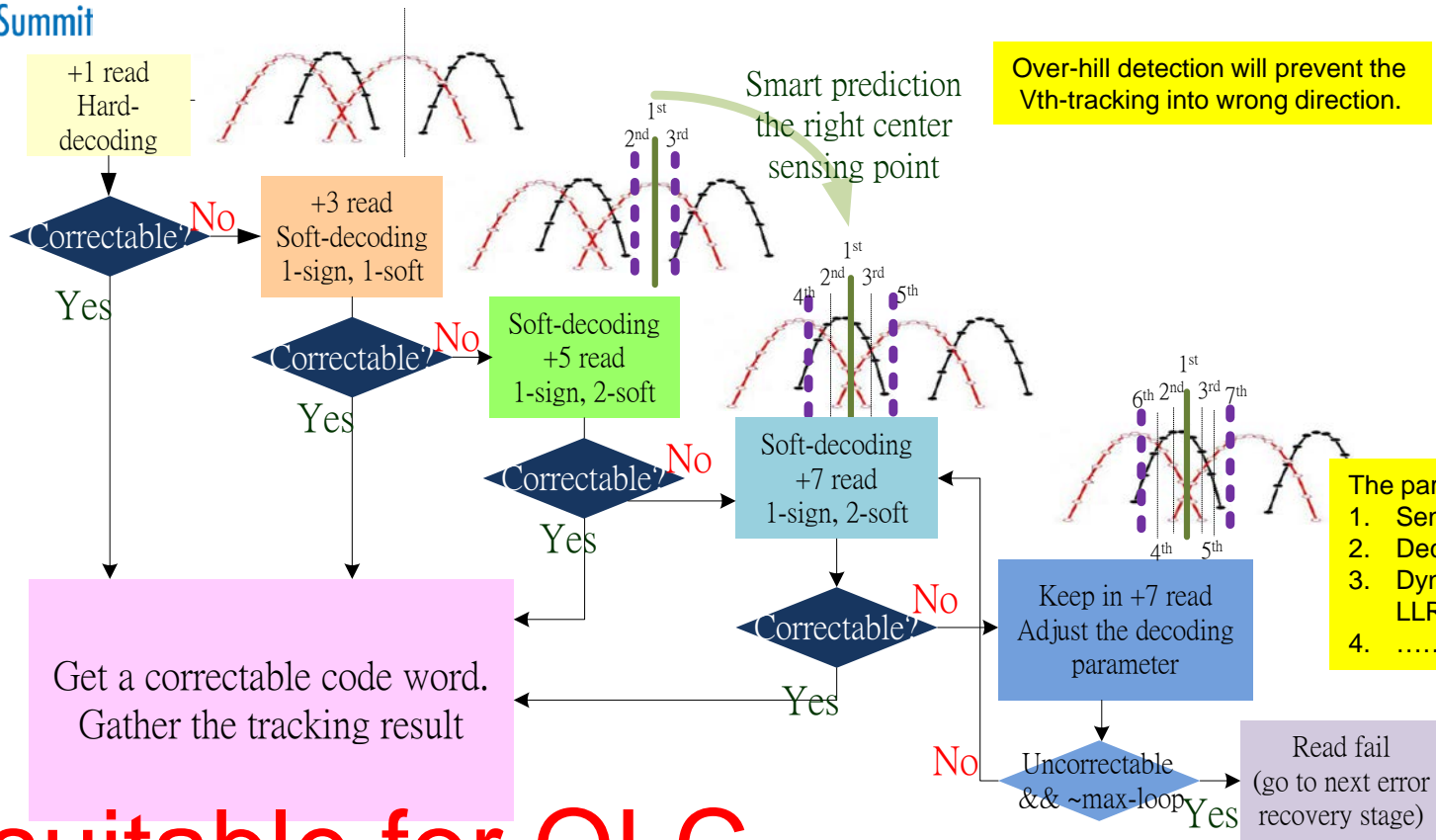


Thread → buffer → NAND





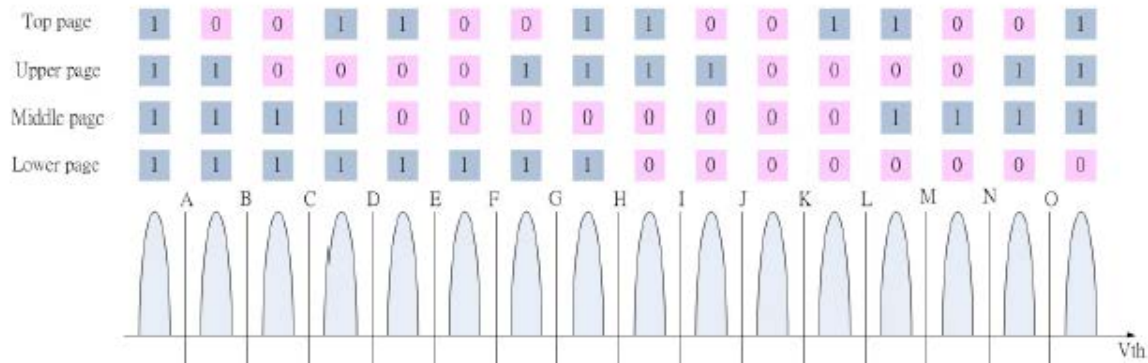
DSP algorithm for the Vth-tracking



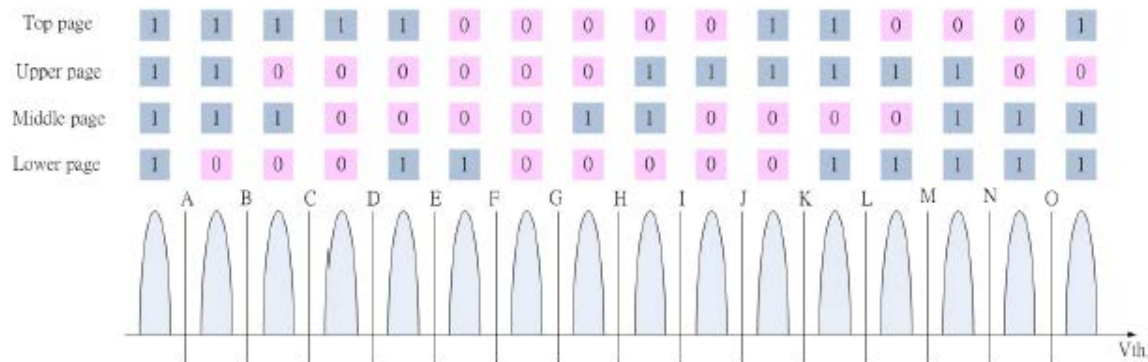
Not suitable for QLC



QLC gray mapping



(a)

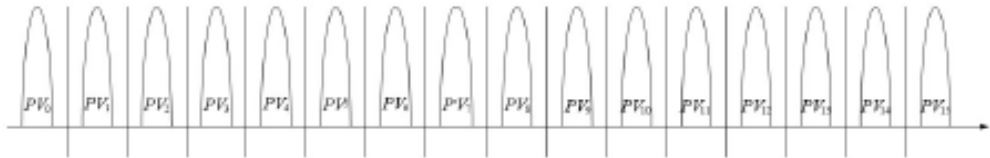


(b)

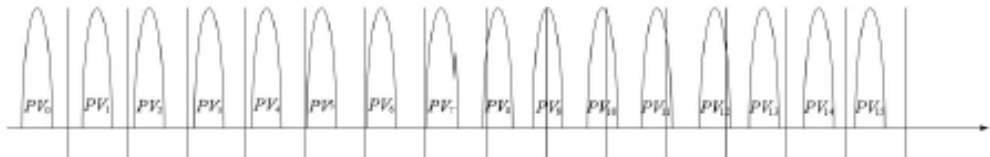


V_{th} distribution shifting case

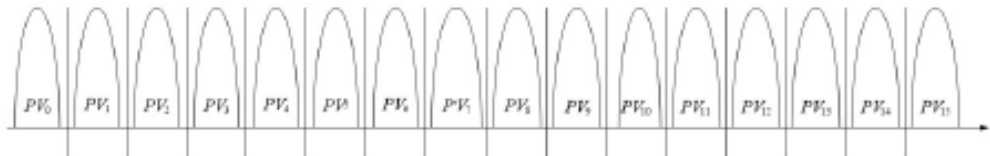
- Traditional V_{th}-tracking is not efficient.



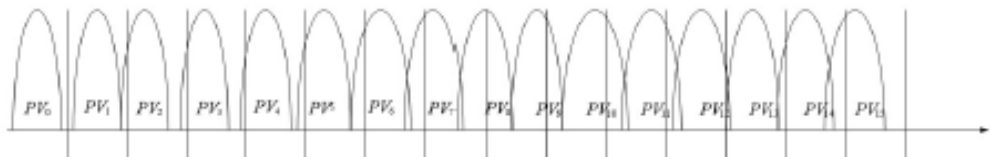
(a)



(b)



(c)



(d)



Controller design future

- Host based or Device based FTL.
- Multi-tenant, guarantee service, latency.
- Always need stronger ECC engine, especially for the hard-info only decoding.
- One sign-bit and one soft-bit still need higher decoding efficiency.
- Advanced process is expensive, single controller should cover three NAND generations



Flash Memory Summit



SiliconMotion

www.siliconmotion.com

Thanks for your attention!
Visit our booth #413 for more information