



Flash Memory Summit

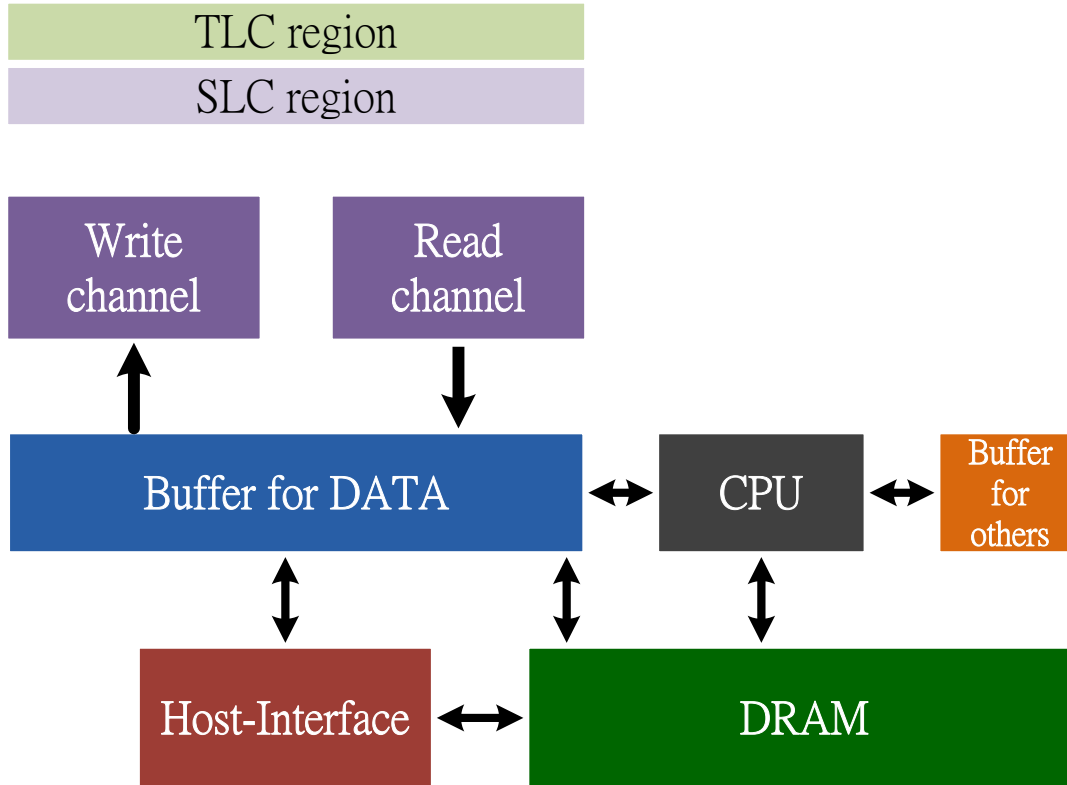


ASIC/Merchant Silicon Chip-Based Flash Controllers

Jeff Yang
Silicon Motion

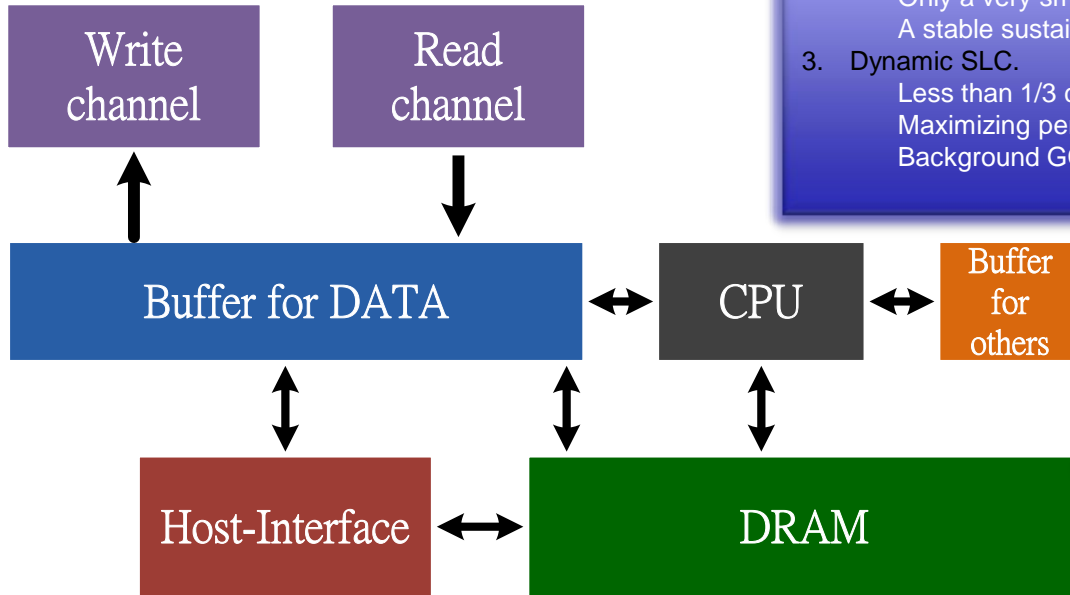


Basic architecture





Application combinations on TLC/QLC



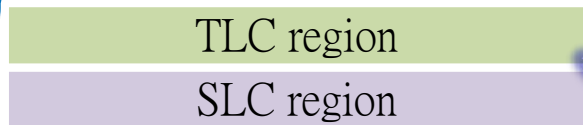
1. SLC caching.
Data always write into SLC. A Fixed portion of SLC regard as a cache buffer.
Performance boost on SLC caching.
Background GC to TLC block.
2. TLC direct.
Only a very small portion for system usage and small random write data.
A stable sustained write performance.
3. Dynamic SLC.
Less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

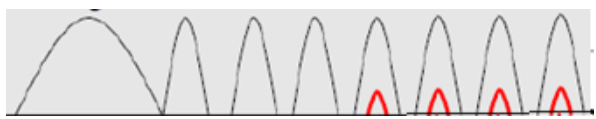
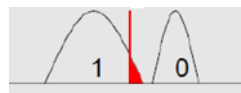
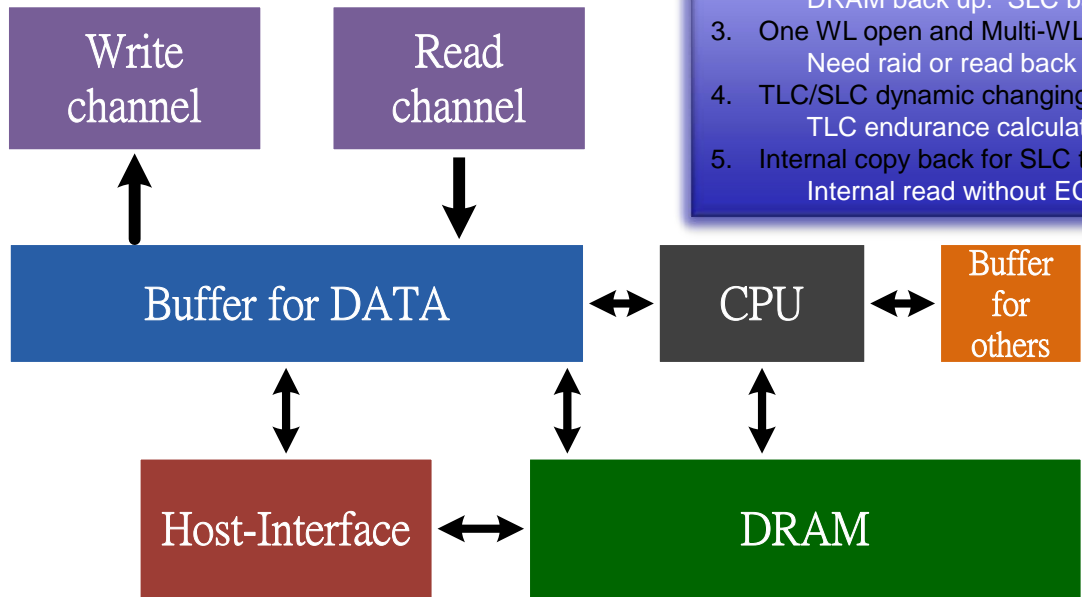


Flash /

Different TLC reliability issue



- 3D.
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





Challenge: Support all combinations and cost efficiency

3D

One-pass

Two-pass

Multi-pass

SLC usage

SLC caching

TLC direct

SLC/TLC
dynamic

External buffer

Full DRAM

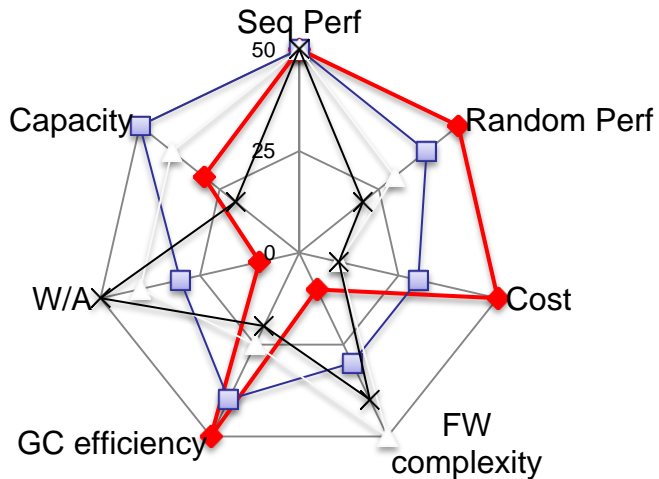
Non-DRAM

Partial DRAM



SMI Product comparison

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

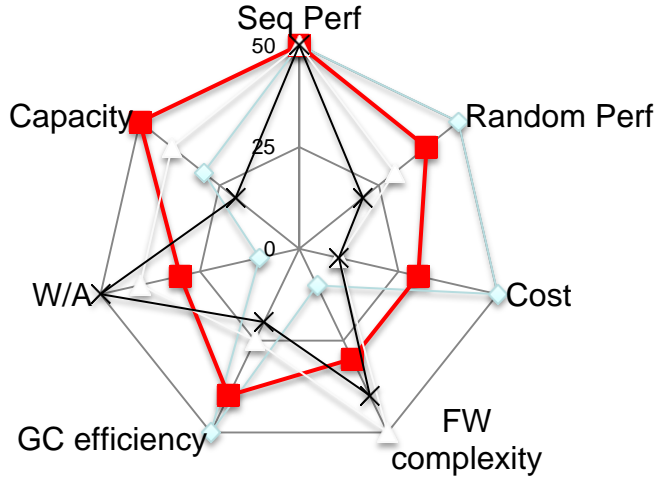


Full DRAM	Partial DRAM	None DRAM	None DRAM (HMB)
Pros <ul style="list-style-type: none">• High Perf• GC efficiency• Low W/A• Low complexity	Pros <ul style="list-style-type: none">• GC efficiency• Capacity limited by L2P address bit (8TB)	Pros <ul style="list-style-type: none">• Cost• Perf still good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



SMI Product comparison

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

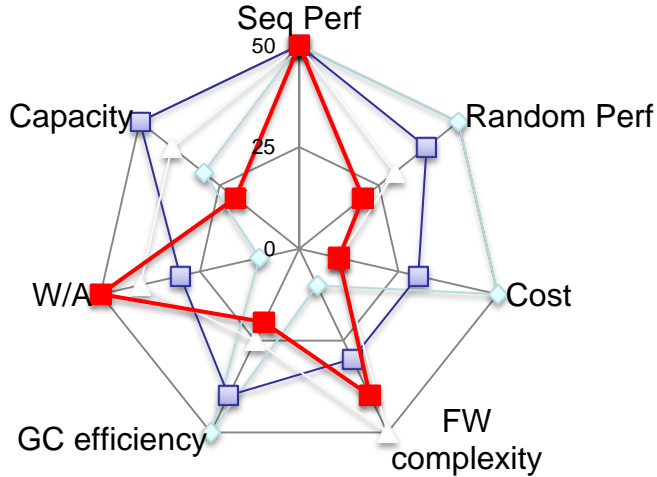


Full DRAM	Partial DRAM	None DRAM	None DRAM (HMB)
<p>Pros</p> <ul style="list-style-type: none"> • High Perf • GC efficiency • Low W/A • Low complexity 	<p>Pros</p> <ul style="list-style-type: none"> • GC efficiency • Capacity limited by L2P address bit (8TB) 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Perf still good for CDM 	<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"



SMI Product comparison

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

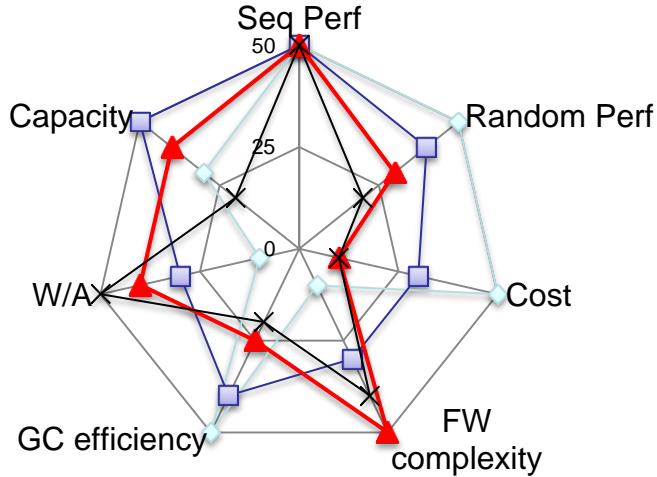


Full DRAM	Partial DRAM	None DRAM	None DRAM (HMB)
<p>Pros</p> <ul style="list-style-type: none"> • High Perf • GC efficiency • Low W/A • Low complexity 	<p>Pros</p> <ul style="list-style-type: none"> • GC efficiency • Capacity limited by L2P address bit (8TB) 	<p>Pros</p> <ul style="list-style-type: none"> • Cost • Perf still good for CDM 	<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"



SMI 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 CDM 	<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"

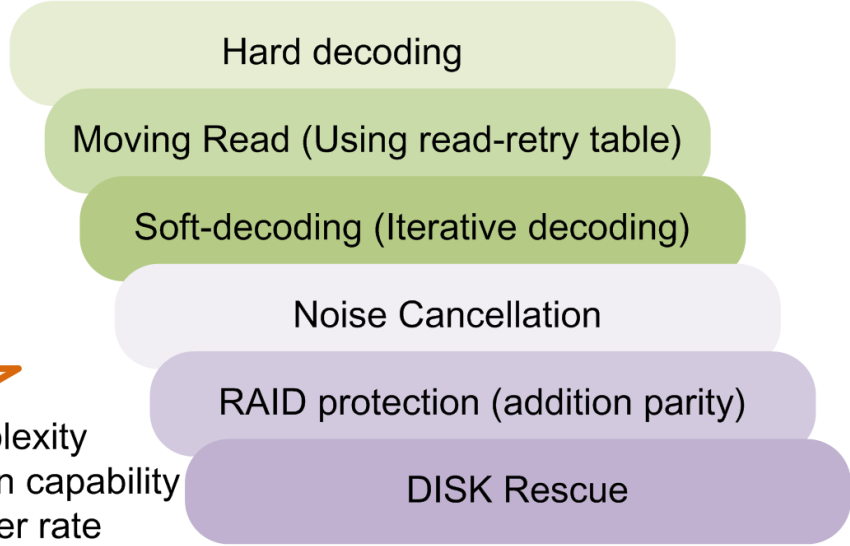
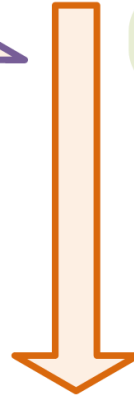
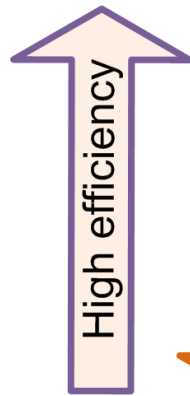


Flash Memory Summit

Traditional Error recovery flow



Error recovery flow



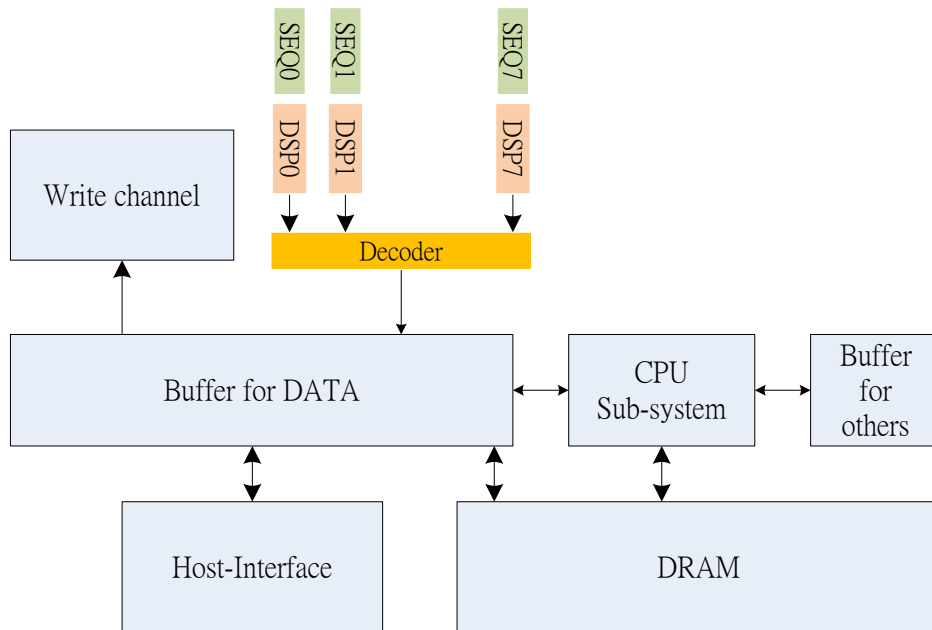
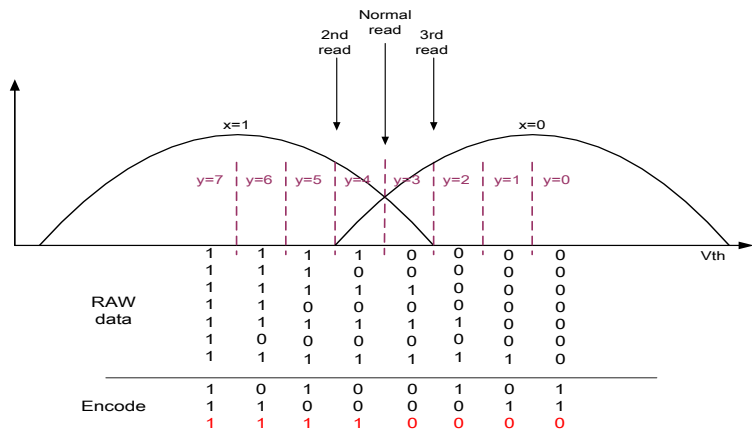
High complexity
 Strong correction capability
 Lower trigger rate

Flash Memory Summit 2013
Santa Clara, CA



Soft-information interface

Flash Memory Summit

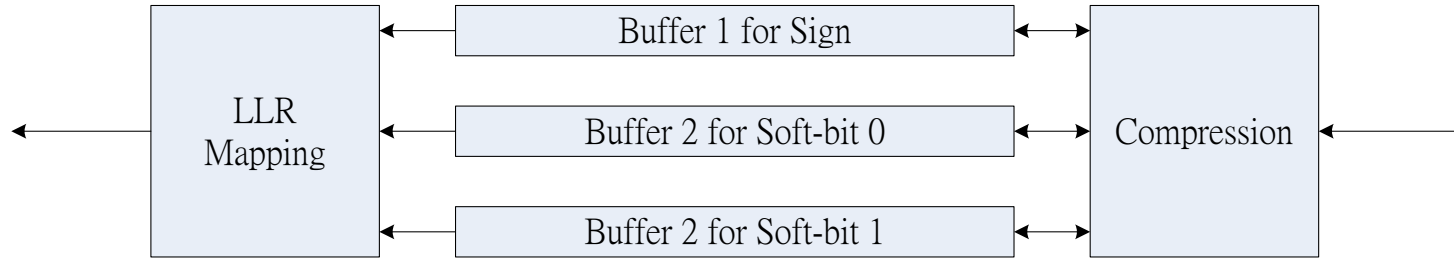


- In order to provide better decoder's correction capability, using the soft-info to get more reliability bits.
- NAND interface support .
 - Traditional read/retry interface.
 - Direct soft-info interface.

DSP engine's buffer size



Flash Memory Summit

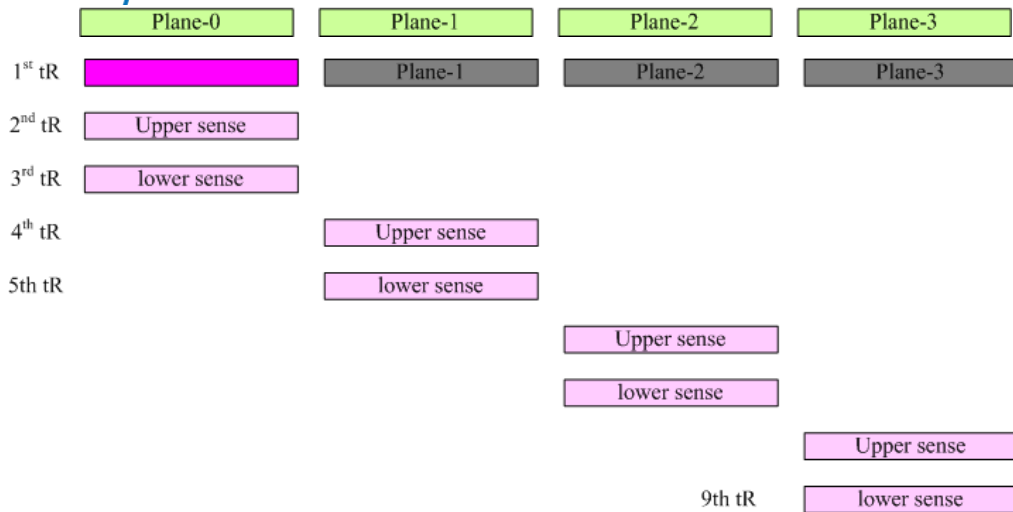


- The buffer size is the capability to contain the number of chunks soft-bit.
- Access addition soft-info from NAND may need additional read busy time.
- Read the soft-bit under the same busy time will have higher efficiency, but buffer size requirement is huge.

Soft-decoding throughput limitation



Flash Memory Summit

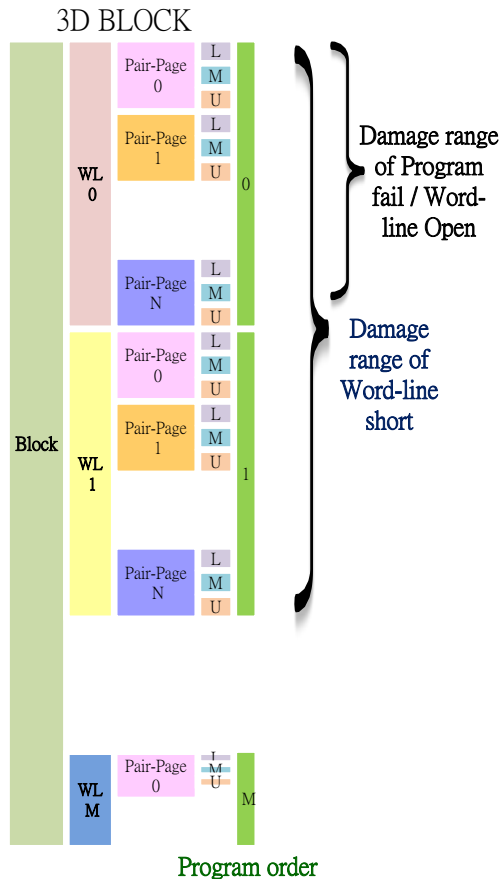
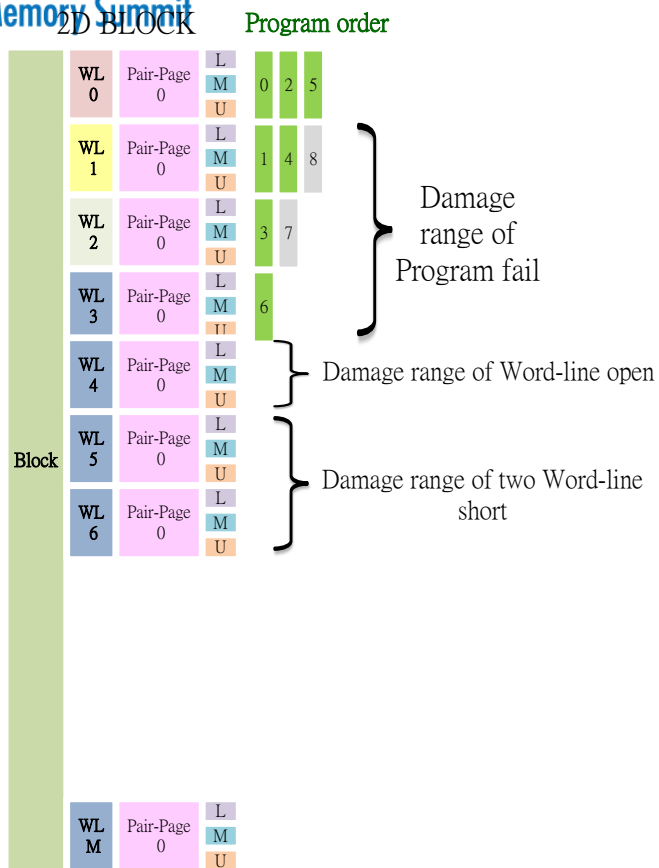


- One Transfer time = $2.5\text{ns}/1\text{B} \times 18432\text{B} = 46\mu\text{s}$ (400MTs)
- Assume DSP-buffer size 16KB.
 - $9 \text{ tR time} + 12 \text{ transfer-time} = 9 \times (100\mu\text{s}) + 12 \times (46\mu\text{s}) = 1452\mu\text{s}$
 - Throughput = $64\text{KB}/1452\mu\text{s} = 44\text{MB}/\text{sec}$
- Assume DSP-buffer size 64KB.
 - $3\text{tR time} + 12 \text{ transfer-time} = 3 \times (100\mu\text{s}) + 8 \times (46\mu\text{s}) = 668\mu\text{s}$.
 - Throughput = $64\text{KB}/668\mu\text{s} = 95\text{MB}/\text{sec}$

In Client SSD applications,
Soft-decoding will regard as the ERROR-Recovery flow.
We will not ask the throughput under recovery mode.
But we will take care the recovery mode trigger rate.



Failure range from 2D to 3D.

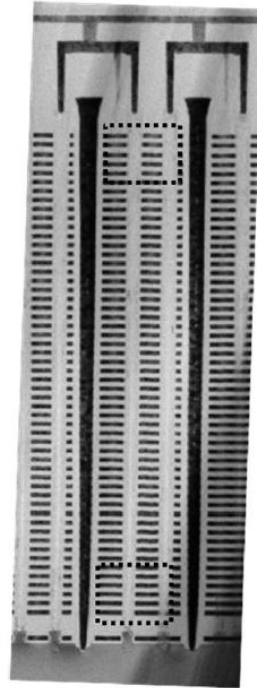
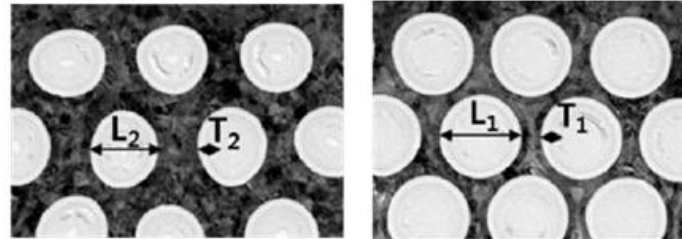




3D NAND Challenges

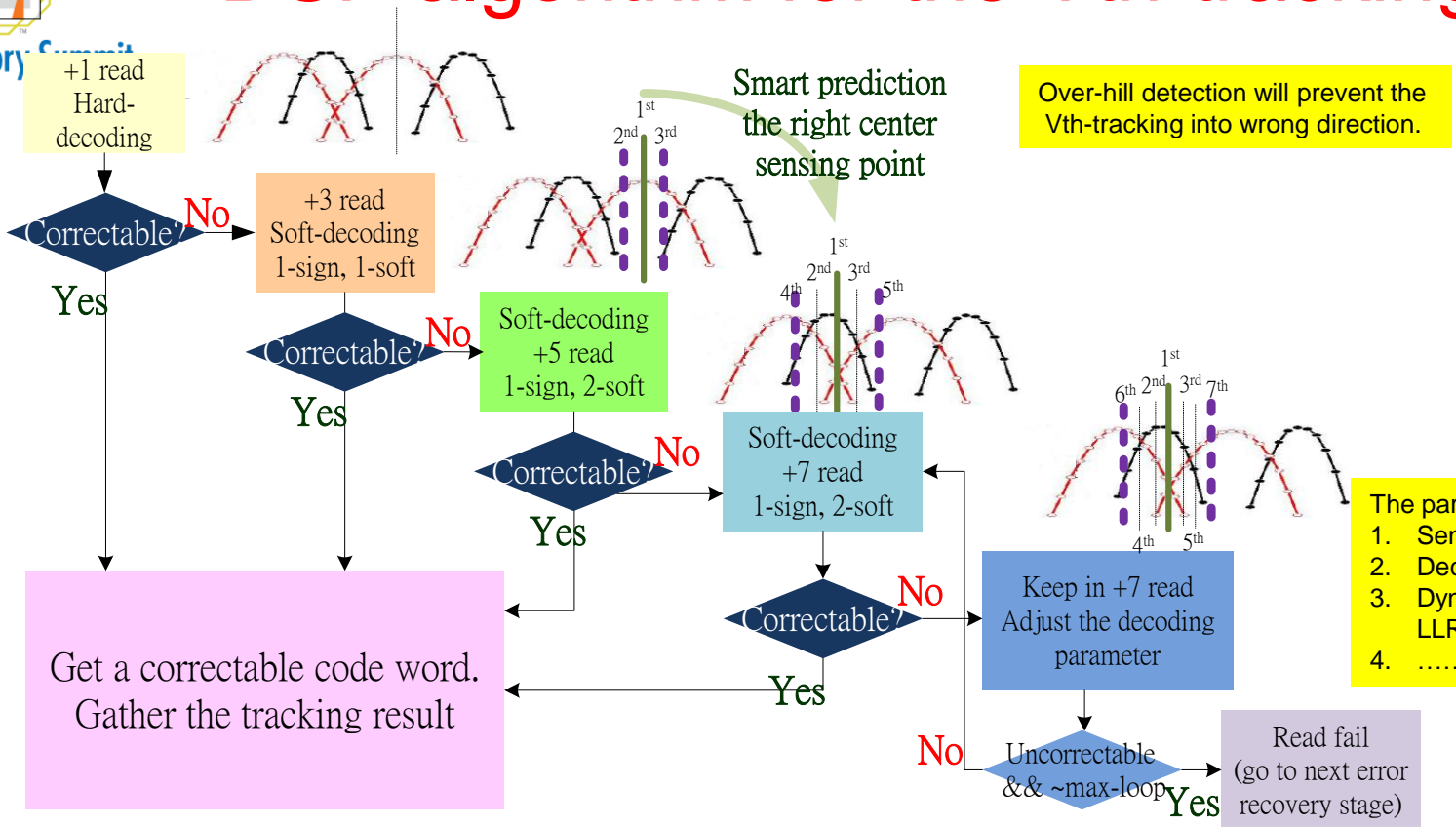
- Each 3D generation will increase the layer number by 30~50%.
- High-aspect ratio channel hole etch.
- Cell current reduction is seriously concerned.
- Reduce the read-voltage to improve the read-count (degradation the read-disturbance) make cell current worse.
- Different cell characteristics for each WL. (program-speed, cell-to-cell interference, retention)
- Poor retention characteristics.

Not easy to screen out some defects.
Especially on bit-column related defect.





DSP algorithm for the Vth-tracking

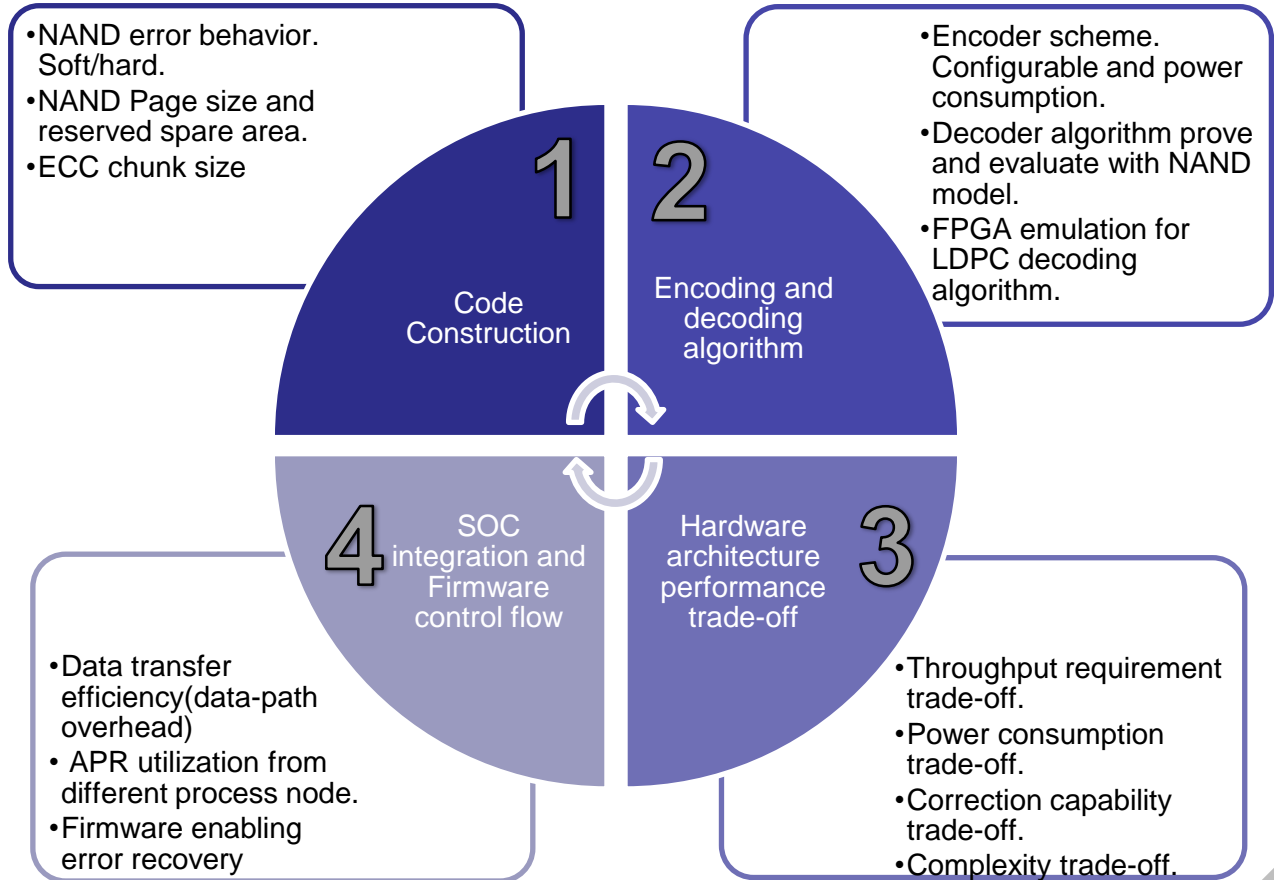




ECC design loop related to NAND characteristics.

Flash Memory Summit

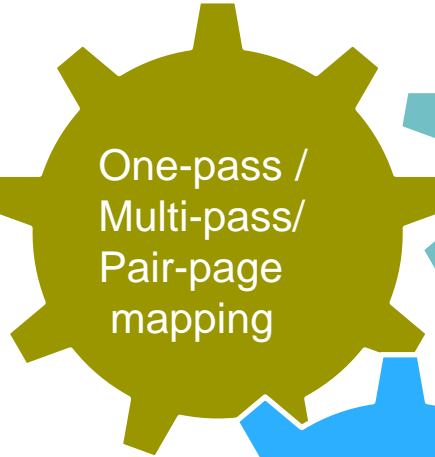
- We already have 7th generation LDPC decoder.
- Keep improving the LDPC performance.
- For higher throughput ~8GB/sec, we may go back to step1.
- After 28nm process, the design iteration depth will from code-construction to trial APR.
- EX: Find the Routing congestion issue in step 4, it may need to solve from step1.



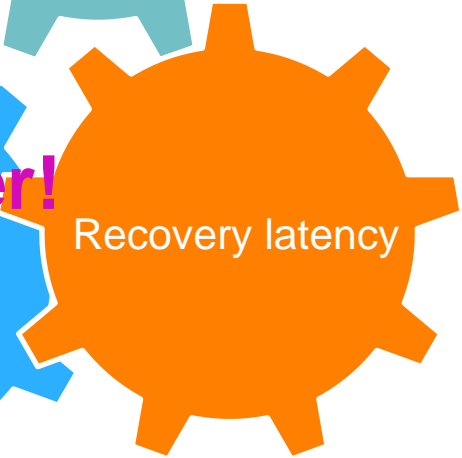
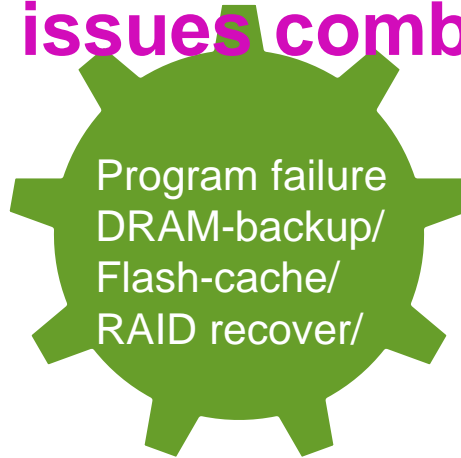
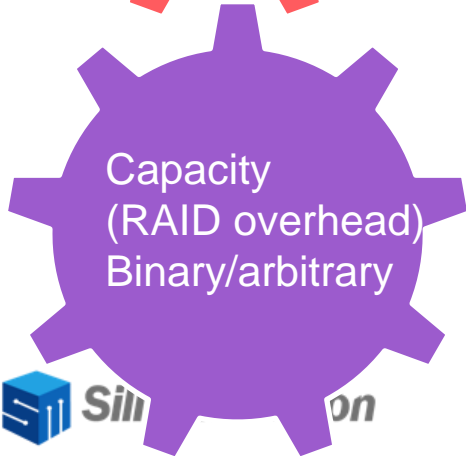
Before the RAID protect flow.....



Flash Memory Summit



All the issues combine together!

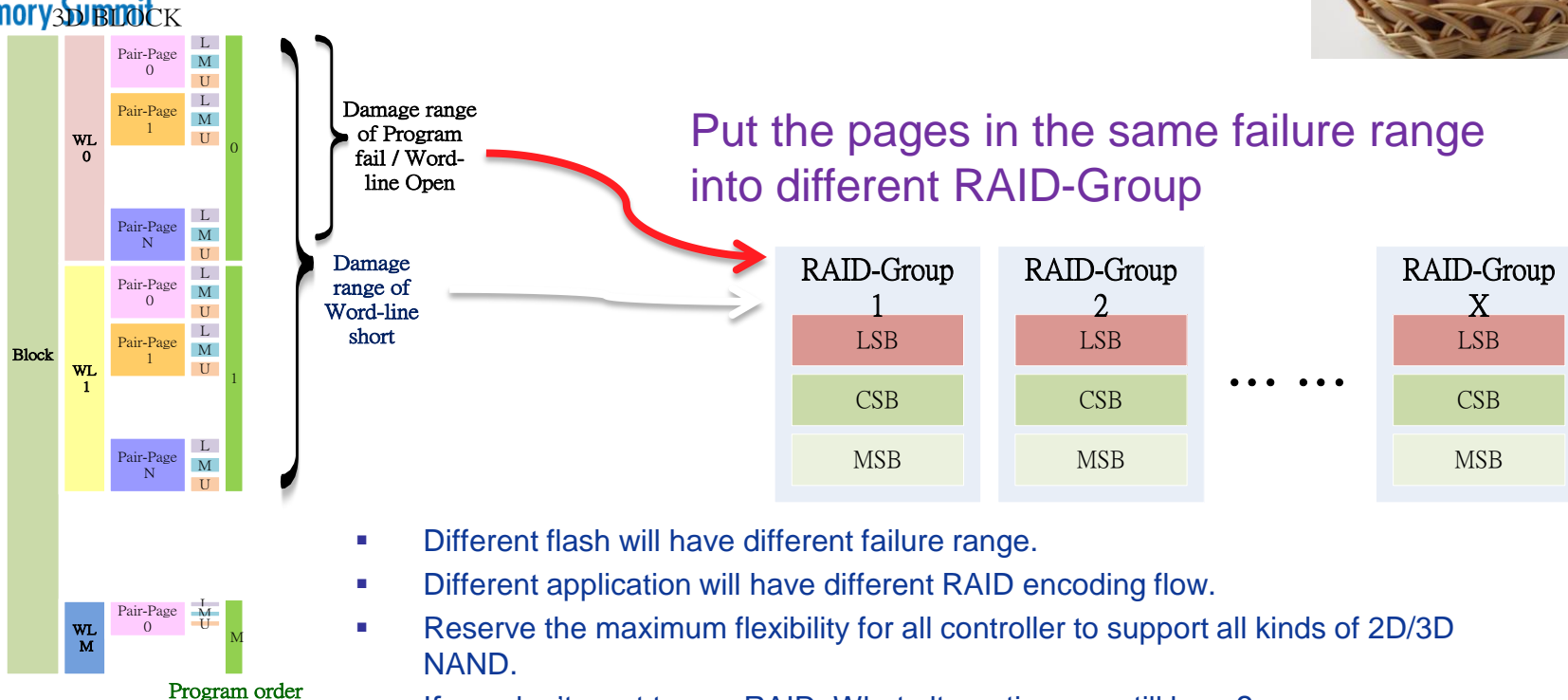


BUT THE SAME CONCEPT IS.....





Don't put eggs in the same basket.



- Different flash will have different failure range.
- Different application will have different RAID encoding flow.
- Reserve the maximum flexibility for all controller to support all kinds of 2D/3D NAND.
- If you don't want to use RAID, What alternative you still have?
 - Read-back check after program.



Matlat demo platform.

Flash Memory Summit

Matlab
interfa
ce

- Matlab base program.
- Gather the data and process by matlab.
- More easily to show the figure.
- Provide high level controller interface.
- NAND access CMD packet

CMD
translati
on

- Hookup the Matlab API to USB/PCIE driver.
- Provide several different NAND access functions to provide Matlab high level control interface.
- SMI controller specification function demo. (Adaptive charge/Vth-tracking/LDPC soft-decoding). Provide a function call.

SATA/P
CIE
controll
er

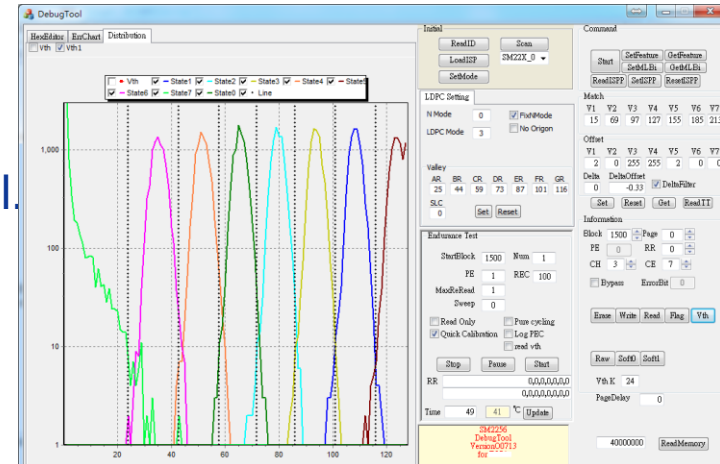
- Firmware implement the Vendor CMD to serve all kinds of NAND access.
- Complex function implement by ARC/ARM firmware to provide specific function demo.



ECC tool update with all NAND.

Flash Memory Summit

- SATA 2258/2259
- PCIE 2260/2262/2263/2270/2264
- Color Vth-distribution. All WL Vth-distribution.
- Error recovery golden flow.
- Soft-info fetch correctness analysis.
- Automation RTBB analysis.
 - (OEM special request support)
 - ECC decoding/encoding, randomizer model.
- ISPP tuning for verification.





Flash Memory Summit

Thanks.

- Q&A