



Flash Memory Summit

# Maximizing Performance of Enterprise QLC-Based SSDs with Flash Processing Units

Jeff Yang

Silicon Motion

- The relation between enterprise SSD and single NAND performance.
- Where is the latency come from.
- A list of the QLC's characteristics on disturbance.
- How to overcome these disturbance.
- Multiple task should be processed in parallel.
- A perfect Flash processing Units provide wonder performance on enterprise SSD.



# Relation Between SSD and NAND

- Read busy effects the SSD read IOPS.
- Program time effects the SSD write IOPS.
- Erase time effects the write IOPS, but much less than Program time.
- NAND IO speed: 2400MTs, 3200MTs, 3600MTs, 4800MTs dominate the IOPS.
- NAND IO CMD protocol overhead impact the IOPS.
- Number of parallelism can be trigger in single busy time increase the IOPS (multi-plane and number of independent CMDs )

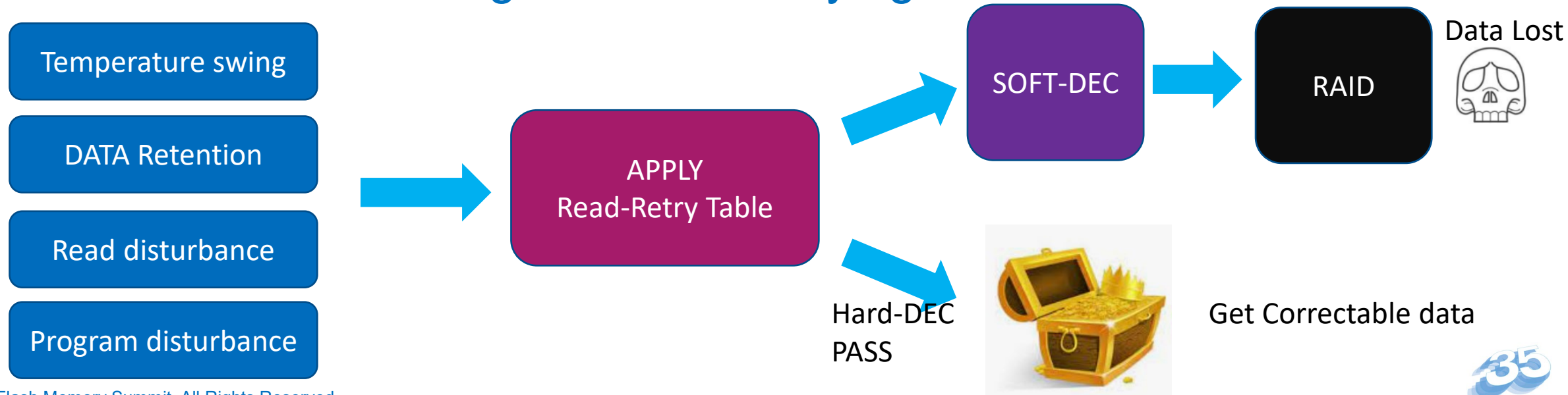


# Where is the Latency come from ?

- Assume the enterprise SSD with DRAM and perfect PLP(Power Lost Protection), the host write data will be buffer in DRAM and return ack immediately.
- $t_{\text{Erase}} > t_{\text{Pro}} \gg t_{\text{R}}$ .
  - During the  $t_{\text{Erase}}$ , and  $t_{\text{Pro}}$ , we want to read the data.
  - During the  $t_{\text{R}}$ , we want to read the data in the same block or plane. (Read collision )
- During the Erase, we can apply the Program. and during the Program we can apply the read. (Suspend feature)
  - Under the Nested Suspend, The CMD interval overhead still dominate the major efficiency.
- During the Read-Collision, the  $t_{\text{R}}$  and the collision rate dominate the main factor.
  - Partition to different physical range on different types of read. (IO-isolation)
  - Reduce the  $t_{\text{R}}$ . (Big problem on QLC)
  - A Smart duplicate the data with high hit rate.

# Latency from Error recovery

- Error bit increase factor.
  - Vth-shifting need better read Voltage from Read-Retry Table,
  - Much smaller Margin between different state need Soft-decoding.
- Error recovery flow is the biggest MONSTER to eat a lot of tR(read busy)
- It will be wonderful if single tR can always get data.





# A list of the QLC's characteristics on disturbance.

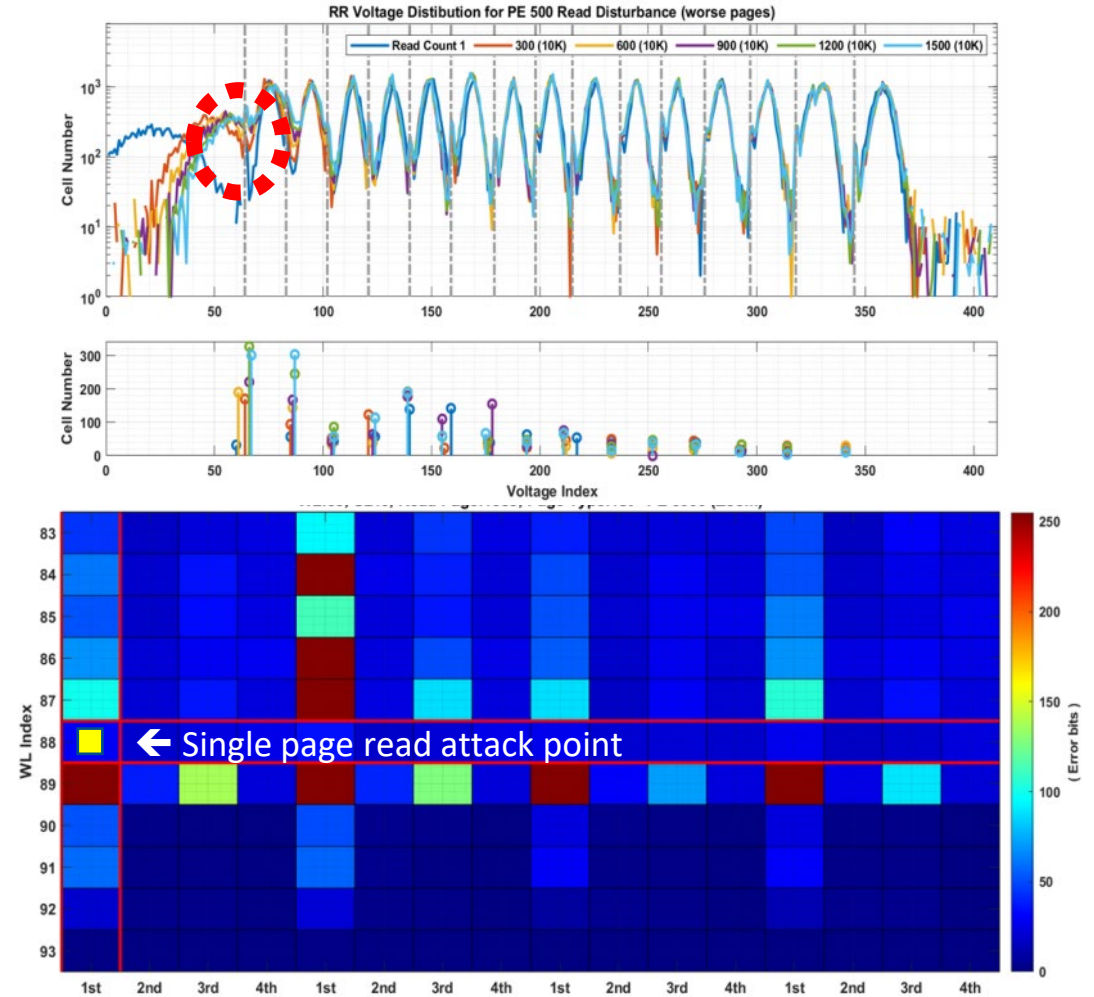
- Multi-pass program mitigate the Program disturbance.
- Read disturbance.
- Room temperature data retention.
- Temperature swing between Program and Read.
- Dynamic SLC/QLC usage for performance boosting.
- A combination between
  - Temp swing cause read/write under different temp 40~70°C
  - Data retention under 40~70°C
  - Read disturbance on the same block or on the same page.

# Read disturbance.



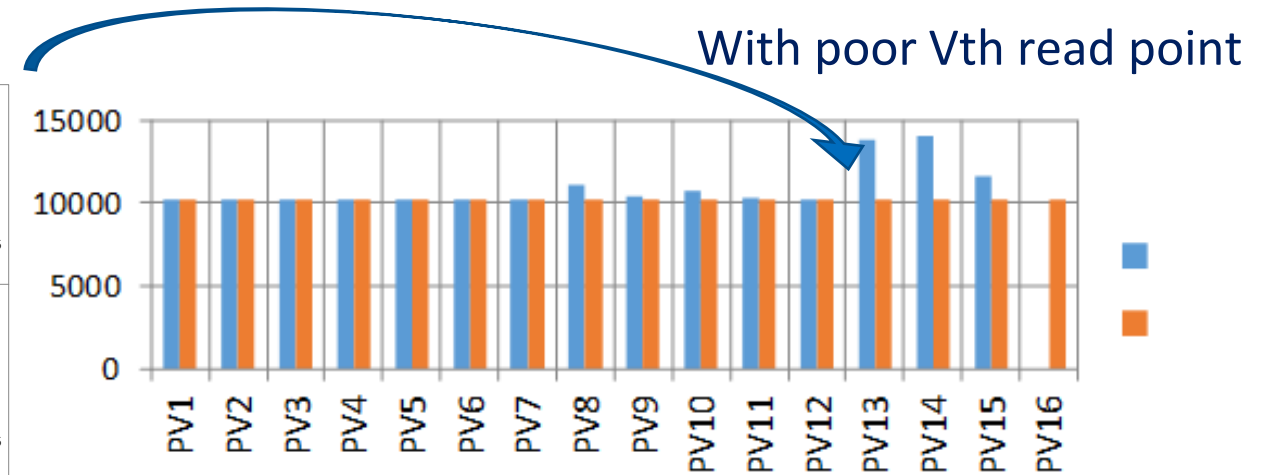
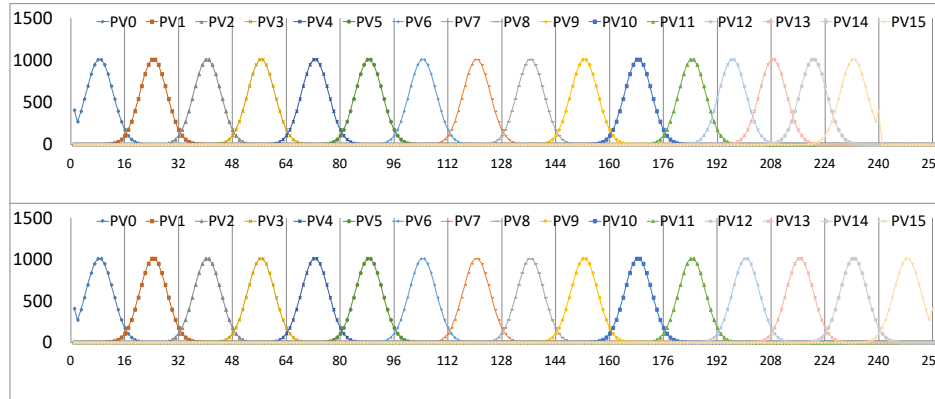
Flash Memory Summit

- QLC will have a large block size. When achieve read-count threshold, move the block also cause big effort.
- A probability based detection algorithm on read disturbance and apply weak neighbor Media scan.
- Only move the data in damage location and mark as a invalid physical address.
- Then, there is no data on the weakest neighbor location. It can tolerant more read attack.



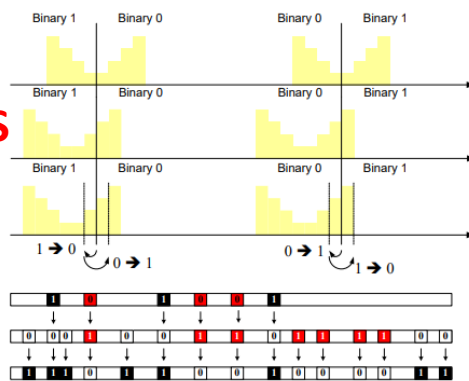
# Finding the Best RV more efficiency.

Data retention



Detecting CSB Direction and Shift Amount

2011 FMS



Too many sensing point in QLC  
tR is much longer in QLC.  
Temperature sensitive.  
More than 100 sets RRT

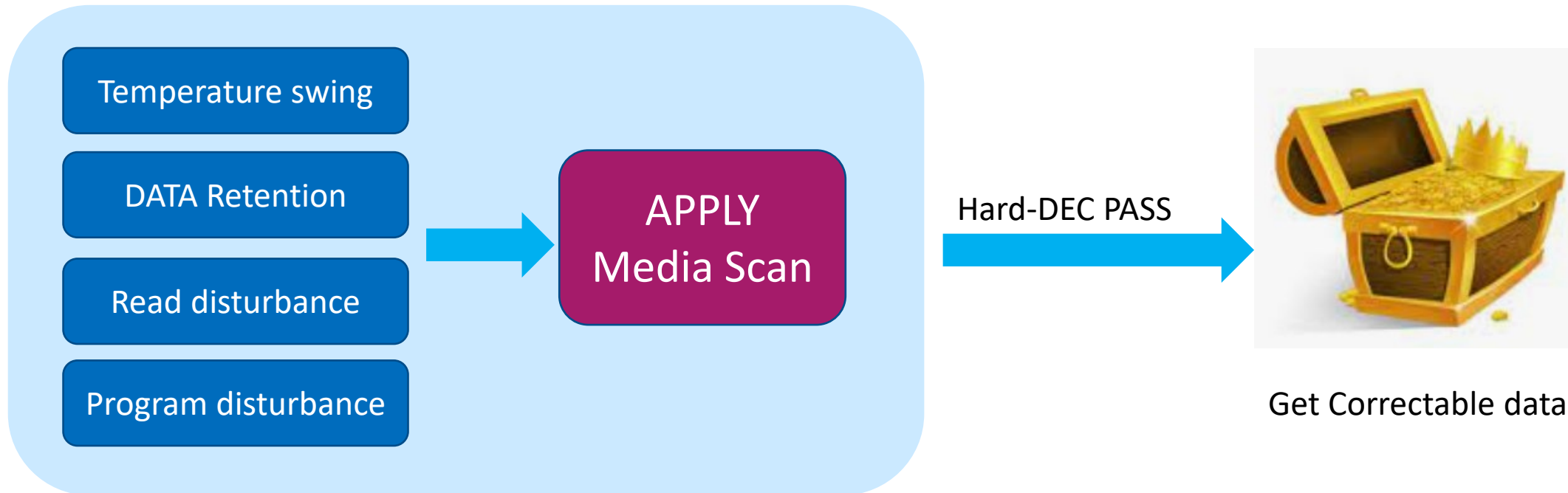
ML  
Algorithm  
With predefine model

Consume MUCH LESS Read



# Media scan makes single-tR to get right data.

- Consume 1% read IOPS resource to provide 99.9999% Hard-DEC PASS



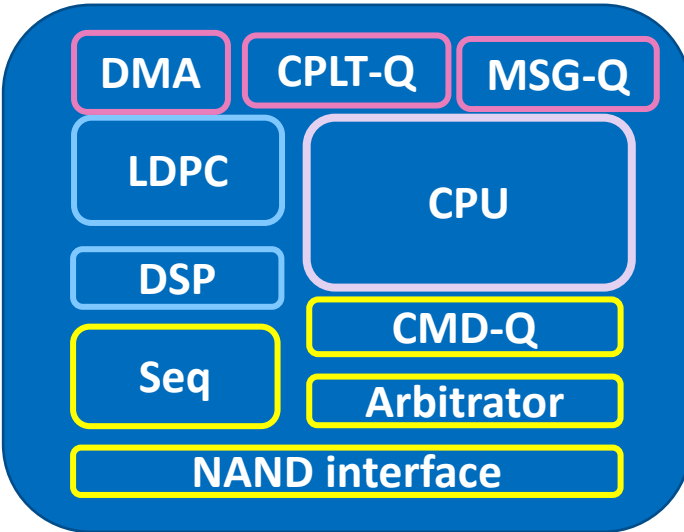


# Multiple task with parallel processing

- Host read (should return data ASAP.)
- Media scan read. (periodical, low weight read)
- GC read. (not urgent, but need maintain a steady traffic )
- Host write. (Need maintain a steady traffic)
- Temperature monitoring on read and write. (depend on temp swing)
- Read disturbance detection read. (depend on read behavior)
- Single Channel with Four physical LUN
  - 6 plane per LUN will have 24 independent task.
  - 4 LUNs with 4 independent Program and Erase task.
  - Each LUN will has its only suspend task. (4 task)



# Flash Processing Unit



Universal interface to main system:

- MSG-Q receive the request.
- CPLT-Q feedback the action status and healthy info of media.
- DMA for DATA move.

CPU:

- Media Algorithm Armed.
- Efficient Media info format for Temp, Retention, Read-count, Endurance.
- Coprocessor for ML engine.

LDPC DSP sub system:

- Perfect interleaving between hard-dec and soft-dec .
- Gather all kinds of Media scan info's RAW data. .
- Zero latency impact when some other read enter deep error recovery.
- No decoding blockage between channels

Optimized NAND interface control:

- Fully utilize NAND IO efficiency.
- Each physical LUN as busy as possible.
- Mitigate NAND and SSD's peak power.
- Support independent plane access.

# Conclusion

- Flash Processing Unit(FPU) is required for the complicated QLC handling for enterprise applications. .
- FPU provides the flexibility and reusable for all kind of Flash application.
- Performance, latency and power become more predictable and higher consistency.
- It makes the single NAND channel become an error free SSD with sequential write.
- Everyone can enjoy the QLC NAND

