F A D U

Storage: The Quiet Engine of Al

Heechul (Fletcher) Chae FADU, Product Planning

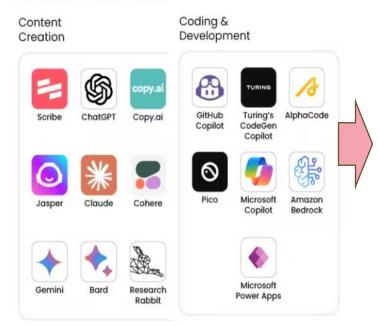


The Shift From Training to Inference

- Generative AI has moved from research to mainstream use, creating billions of inference events daily.
- By 2030, ~75% of Al compute will be inference

— design data centers for high-volume inference traffic.

Generative AI Tools





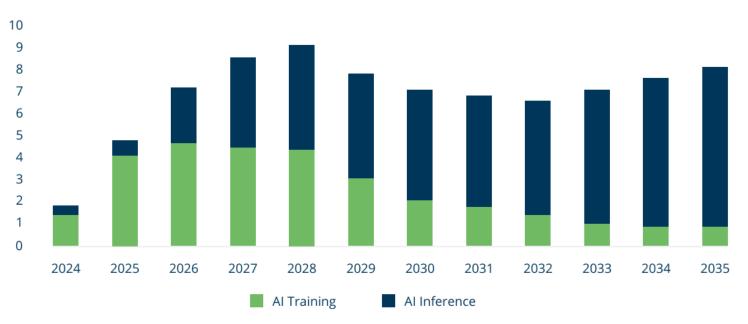


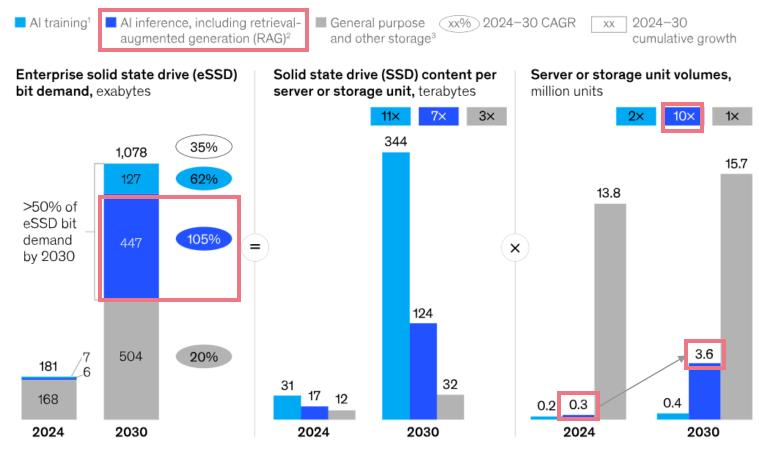
Figure 10: Approximately 75% of Future AI Compute Demand to Come From Inference by 2030

Source: Brookfield internal research.

**Brookfield Building the Backbone of Al.pdf

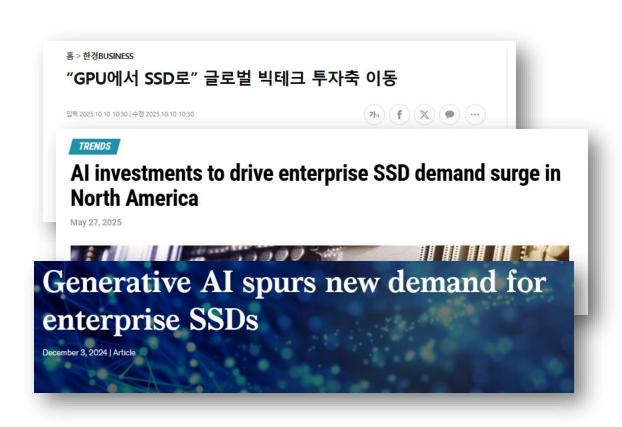
Inference Server (w RAG) is Driving Explosive Data Demand

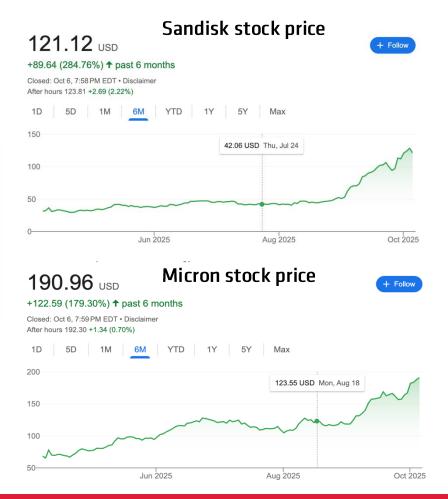
- Al inference infra(GenAl) data demand will grow ~105% CAGR through 2030—over 3.5× training.
- Main driver: RAG—via vector DBs (indexing, replication, low-latency fetch).



Capital Flows Toward eSSDs in the GenAl Supercycle

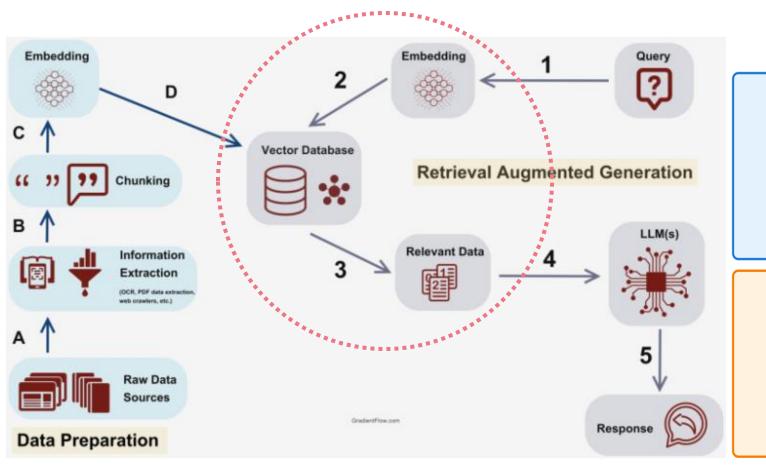
- Investment and attention are converging on eSSDs for the inference era.
- It signals a need for specialized eSSDs—not bulk HDDs or commodity drives—like HBM for training.





RAG Makes Storage the Bottleneck in Inference (1)

- Vector search with query embeddings is crucial for inference QoS.
- To make vector search faster, caching key pieces in HBM and local SSD.



Vector DB Caching Effects

Tail Latency Suppression

Eliminate PCIe/DRAM round-trips
Remove CPU scheduling jitter

→ p99/p999 Stabilized

QPS Increase ↑

GPU massive parallel search (1000s~10000s queries batched)

→ No memory bottleneck

Recall Cost Mitigation

Increase nprobe (IVF)
efSearch (HNSW)

→ Gentle latency growth

Data Copy Cost ↓

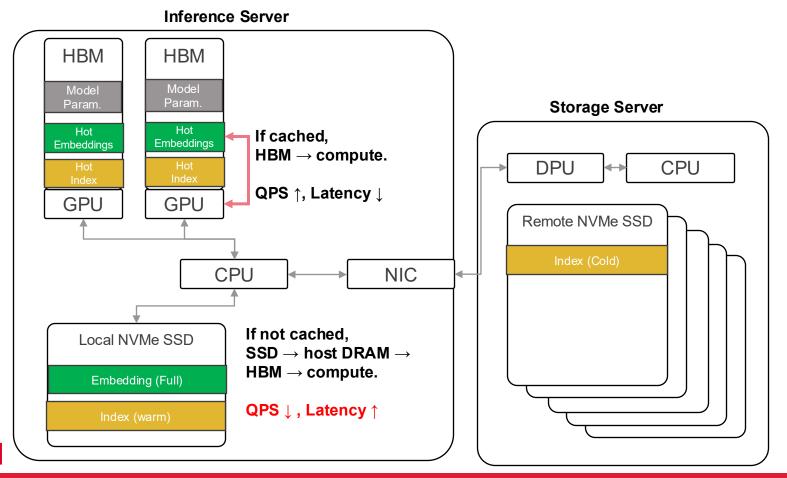
Candidate \rightarrow Distance \rightarrow Re-rank

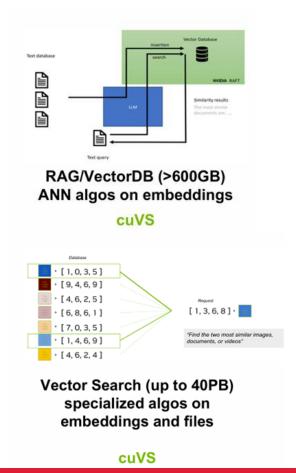
All in HBM pipeline

→ Minimize memory bounce

RAG Makes Storage the Bottleneck in Inference (2)

- As the vector DB grows, caching gets complex causing OOM(Out-of-Memory)
- As I/O shifts from HBM↔GPU to NVMe SSD, SSD BW becomes the bottleneck
 - this is why storage is critical for inference.





RAG Workloads Storage Must Face

Volume

- Provide O(100M) IOPs!
 - Gen5 PCle (x16 lane): 50 GB/s
 - IOPS per GPU
 - = PCIe bandwidth / xfer xsize
 - = 50 GB/s / 512B
 - = 100 MIOPs/GPU
- 예) Blackwell 기준 (2 GPU to 4 NVMe) IOPS per GPU = 200 MIOPS / GPU = 100 MIOPS / NVMe (Can be required)

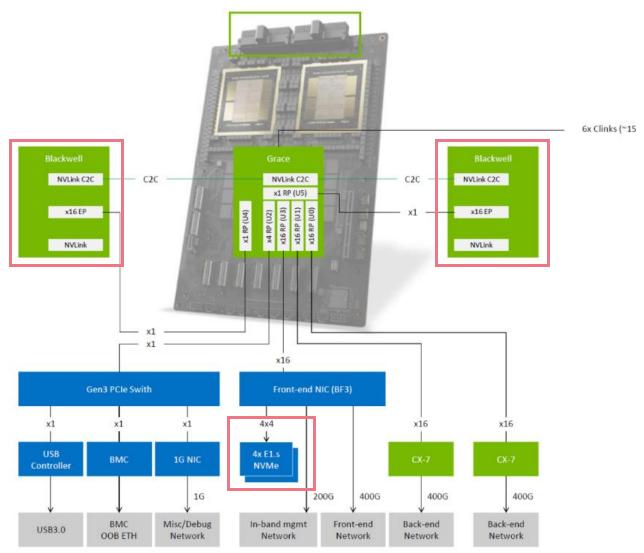
Size

- Under 512Bytes
- Embeddings 512B-4KB, graph structure could be 8B-128B

Pattern

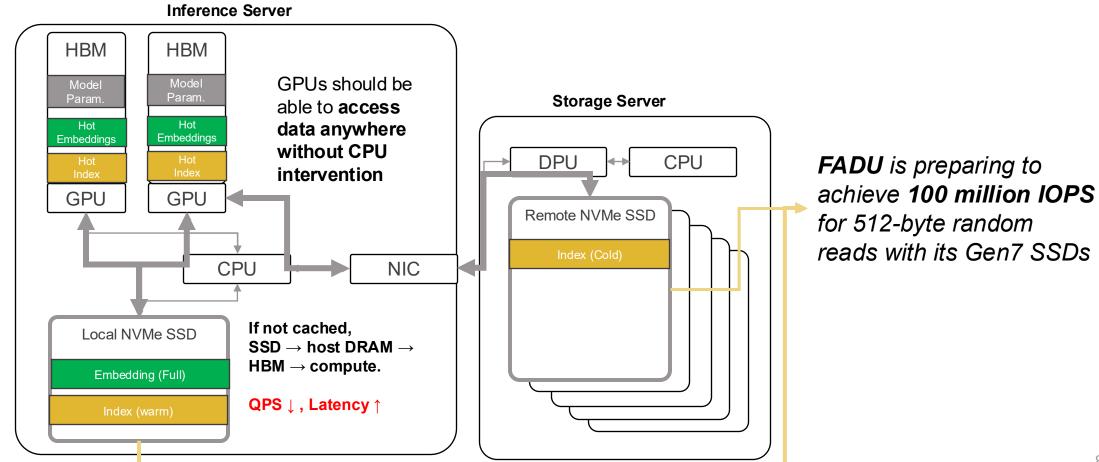
- Sparse (Random)

GB200/300 Compute Tray (View only 1P (2GPU/1CPU)

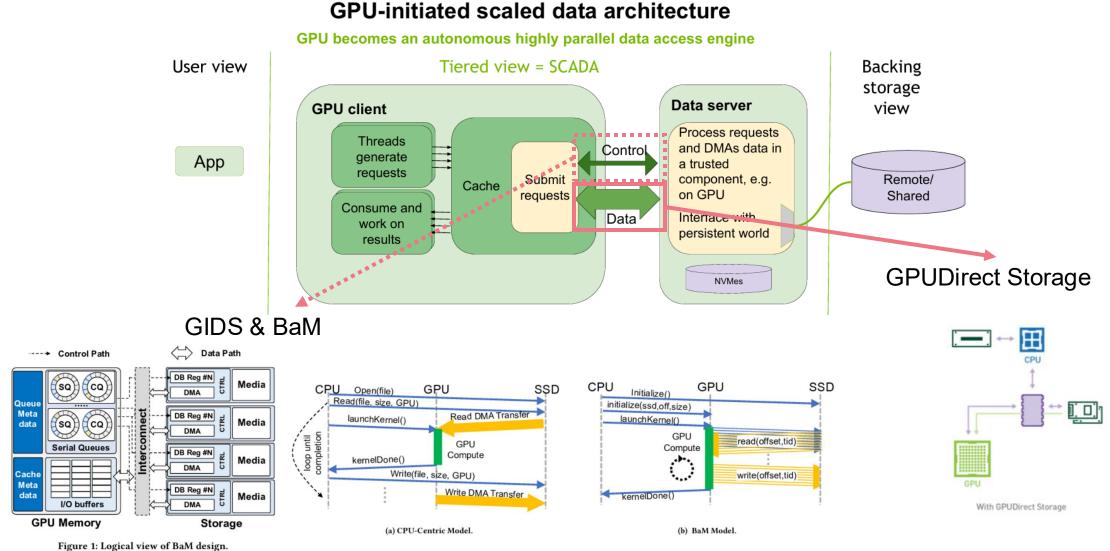


What to solve?

- 1. (Infrastructure Developers) Optimize the path from GPUs to NVMe SSDs.
- 2. (Storage Developers) Make SSDs can handle the I/O rate generated by the GPU.



Related work – Optimizing GPU-to-NVMe SSD Paths



Future works

1. Remove local SSDs from the inference server

- Will the network (e.g., the NIC) become the bottleneck?
- Do we need a new interface to the remote storage server?

2. A complete implementation of GPU-initiated storage access

- GPU-side NVMe command-level access is possible;
- However, without filesystem-level direct access, aren't we still dependent on the CPU?

F A D U

Random Read Maximization for AI Storage

Tech to overcome NAND Channel Bottleneck

EuiJin.Kim, FADU Firmware

Contents

- Random Read IOPS Trend for AI
- Target Random Read IOPS for Future Solution
- To achieve 100 Million IOPS of Random Read
 - Required System Configurations
 - Issue with NAND Channel
- Background: NAND Read Operation Flow
- Why the NAND Channel is the Bottleneck
- NAND Channel Optimization
 - Data I/F Optimization
 - CMD I/F Optimization
- Conclusion

Random Read IOPS Trend for AI Storage

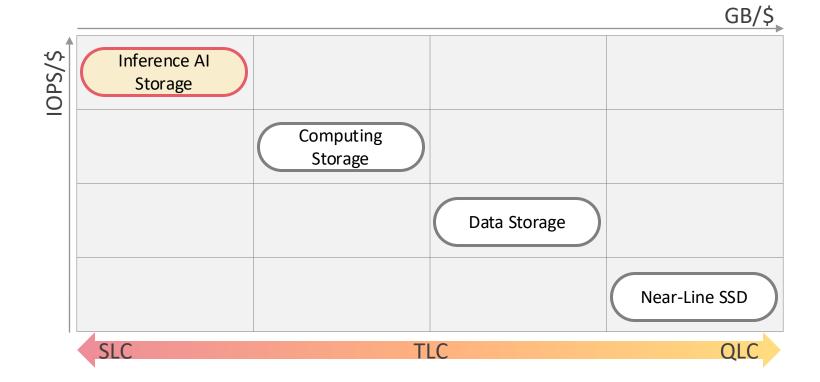
Demand for Inference AI Storage

- IOPS / \$
- IOPS / Watt
- Fine-Grained Random Read Intensive
 - embeddings 512B~4KB
 - graph structure could be 8B~128B

SSD Tech Trend for Inference AI Storage

Focused on "Random Read IOPS" & "Watt"

- 512Byte Block Size
 - More IOPS via fixed speed I/F
- SLC NAND
 - Lowest Cell Read Time
 - Lowest Cell Read Power



Target Random Read IOPS for Future Solution

Maximized BW & IOPS on PCIe Gen7

Block size	Ideal Read BW	Ideal Read IOPS
4KB	55.7 GB/s	13.6M IOPS
512Byte	52.7 GB/s	102.9M IOPS (7.5)

Target Random Read IOPS for Future Solution

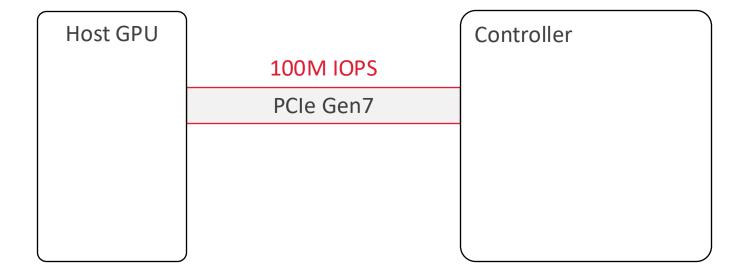
Maximized BW & IOPS on PCIe Gen7

Block size	Ideal Read BW	Ideal Read IOPS
4KB	55.7 GB/s	13.6M IOPS
512Byte	52.7 GB/s	102.9M IOPS (7.5)

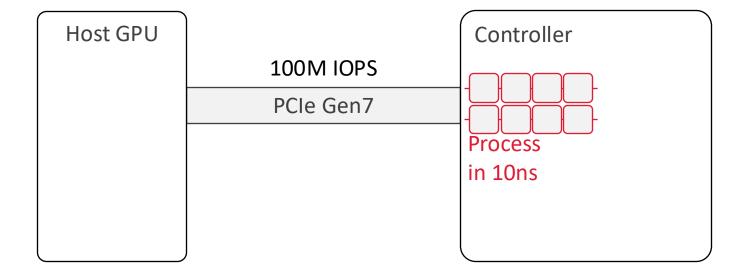
Target IOPS for Future Solution
100M IOPS!!

System capability to maximize Random Read IOPS

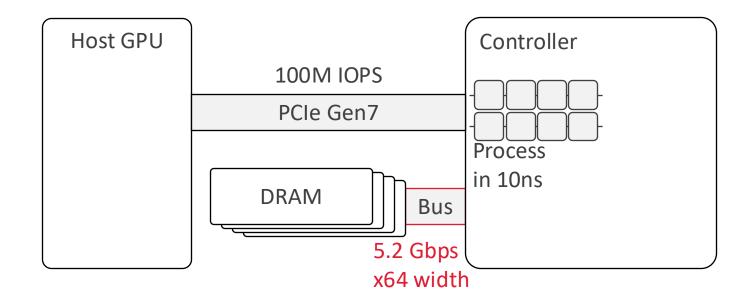
• Fully utilized PCIe Gen7 throughput



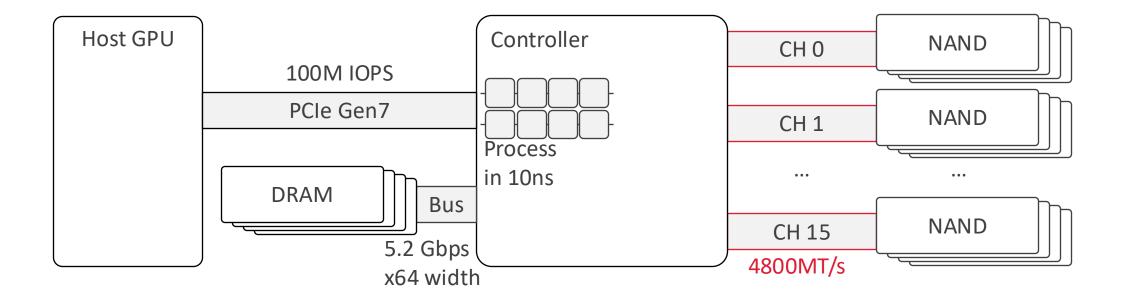
- Fully utilized PCle Gen7 throughput
- Controller processing capability to handle single I/O in 10ns



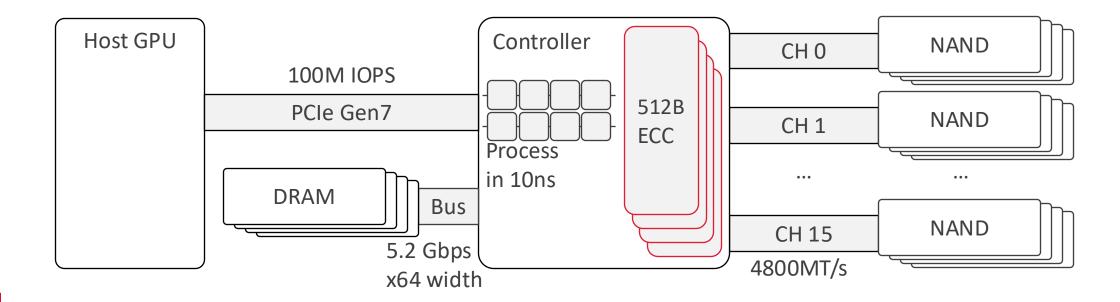
- Fully utilized PCle Gen7 throughput
- Controller processing capability to handle single I/O in 10ns
- DRAM bus throughput for 100M metadata access



- Fully utilized PCle Gen7 throughput
- Controller processing capability to handle single I/O in 10ns
- DRAM bus throughput for 100M metadata access
- Practical solution with 16 NAND Channels x 8 Dies which support 4800MT/s



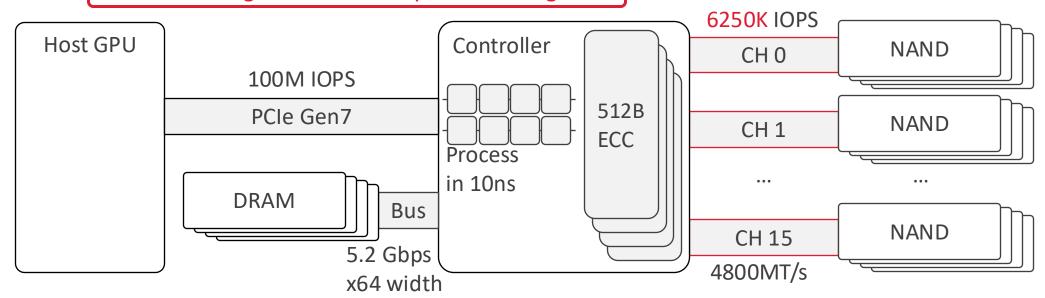
- Fully utilized PCle Gen7 throughput
- Controller processing capability to handle single I/O in 10ns
- DRAM bus throughput for 100M metadata access
- Practical solution with 16 NAND Channels x 8 Dies which support 4800MT/s
- 512Byte ECC Engine Customized for AI workload



System capability to maximize Random Read IOPS

- Fully utilized PCIe Gen7 throughput
- Controller processing capability to handle single I/O in 10ns
- DRAM bus throughput for 100M metadata access
- Practical solution with 16 NAND Channels x 8 Dies which support 4800MT/s
- 512Byte ECC Engine Customized for AI workload
- 6.25M IOPS per Channel (= 100M IOPS / 16 Channels)
 - Required effective BW of single channel: 3600MT/s
 - Required channel efficiency: 75%

It is hurdle to get a CH efficiency of 75% or higher!!



System capability to maximize Random Read IOPS

- Fully utilized PCIe Gen7 throughput
- Controller processing capability to handle single I/O in 10ns
- DRAM bus throughput for 100M metadata access
- Practical solution with 16 NAND Channels x 8 Dies which support 4800MT/s
- 512Byte ECC Engine Customized for AI workload
- 6.25M IOPS per Channel (= 100M IOPS / 16 Channels)
 - Required effective BW of single channel: 3600MT/s
 - Required channel efficiency: 75%

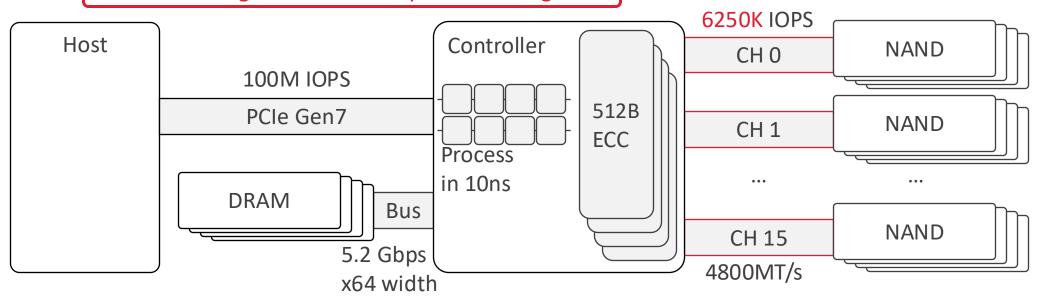
It is hurdle to get a CH efficiency of 75% or higher!!

The easiest way to achieve throughput is to "increase the # of channels"

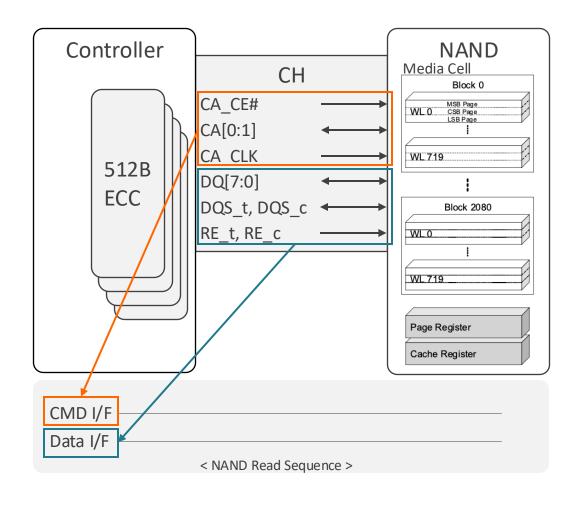
The more channels,

- the higher the throughput can be achieved
- the larger the chip size
- the less installation space

→ Find a way to solve the issue, while practically staying within 16 channels.

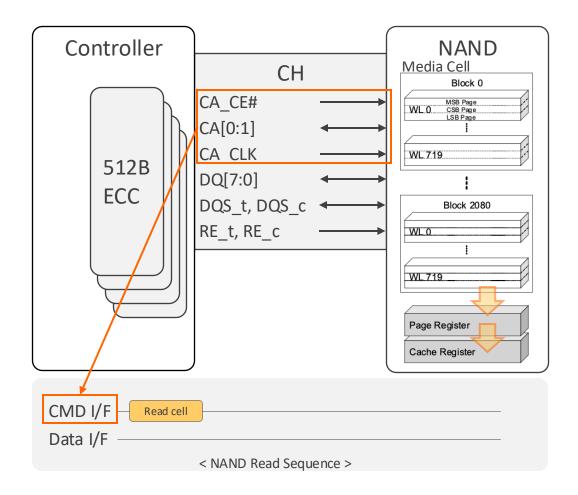


How NAND Read works & Interact with controller (on SCA Protocol)



How NAND Read works & Interact with controller (on SCA Protocol)

- 1. Read Cell Command via CMD I/F
 - Controller designates the cell address to read
 - NAND status: changed to busy
 - Sensing Read data from Cell
 - Load data to page register
 - Copy data from page register to cache register
 - NAND status: changed to ready after finishing job

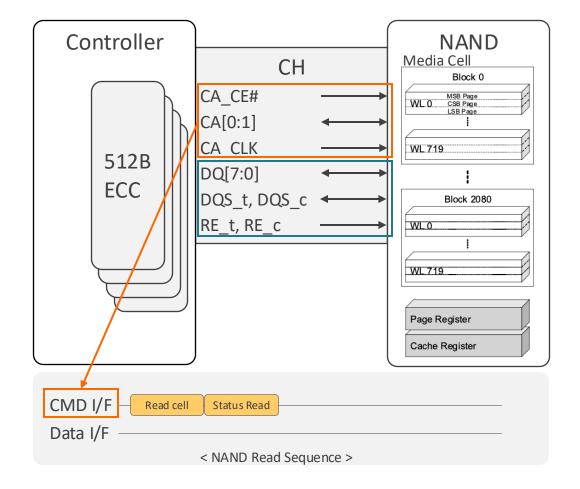


How NAND Read works & Interact with controller (on SCA Protocol)

- 1. Read Cell Command via CMD I/F
 - Controller designates the cell address to read
 - NAND status: changed to busy
 - Sensing Read data from Cell
 - Load data to page register
 - Copy data from page register to cache register
 - NAND status: changed to ready after finishing job

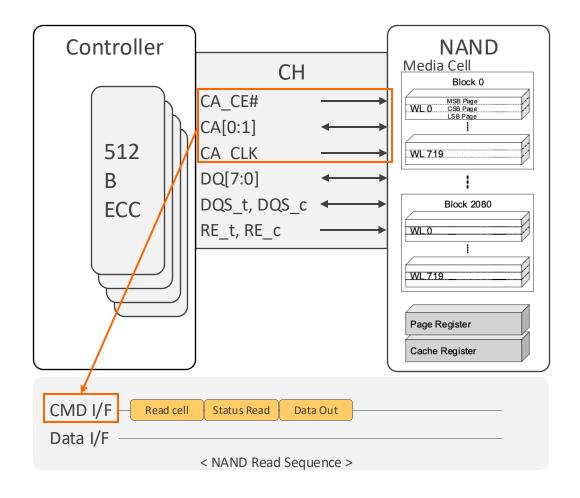
2. Status Read Command via CMD I/F

- Controller check NAND is ready to read-out data
- NAND returns current status



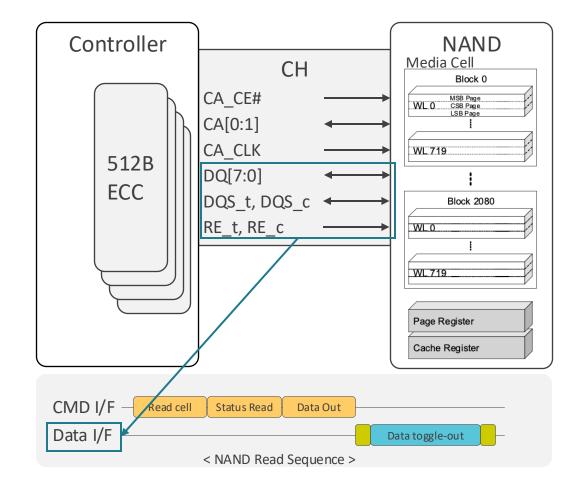
How NAND Read works & Interact with controller (on SCA Protocol)

- 1. Read Cell Command via CMD I/F
 - Controller designates the cell address to read
 - NAND status: changed to busy
 - Sensing Read data from Cell
 - Load data to page register
 - Copy data from page register to cache register
 - NAND status: changed to ready after finishing job
- 2. Status Read Command via CMD I/F
 - Controller check NAND is ready to read-out data
 - NAND returns current status
- 3. Random data-out command via CMD I/F
 - The controller notifies the NAND to prepare to toggle out data.



How NAND Read works & Interact with controller (on SCA Protocol)

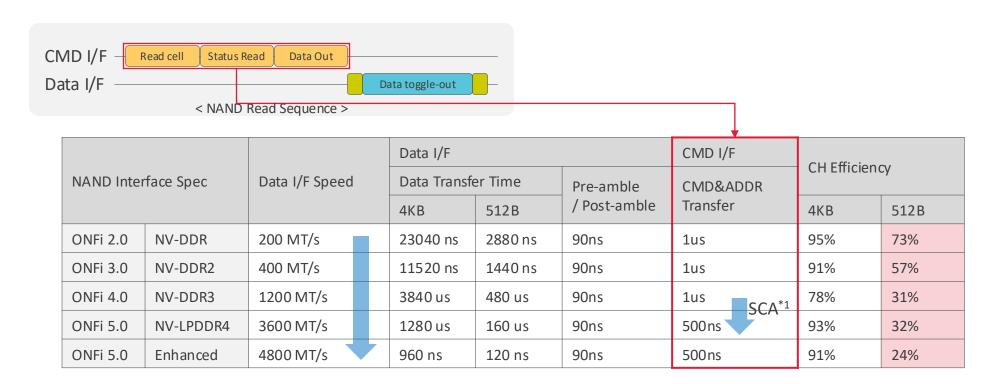
- 1. Read Cell Command via CMD I/F
 - Controller designates the cell address to read
 - NAND status: changed to busy
 - Sensing Read data from Cell
 - Load data to page register
 - Copy data from page register to cache register
 - NAND status: changed to ready after finishing job
- 2. Status Read Command via CMD I/F
 - Controller check NAND is ready to read-out data
 - NAND returns current status
- 3. Random data-out command *via CMD I/F*
 - The controller notifies the NAND to prepare to toggle out data.
- 4. Data Toggle Out via Data I/F
 - Controller toggles RE to NAND
 - NAND toggle-out data in cache register via NVDDR I/F (DQ[7:0])



Why the NAND Channel is the Bottleneck

Factors that reduce channel efficiency

 Slow improvement of the CMD I/F speed with legacy manner

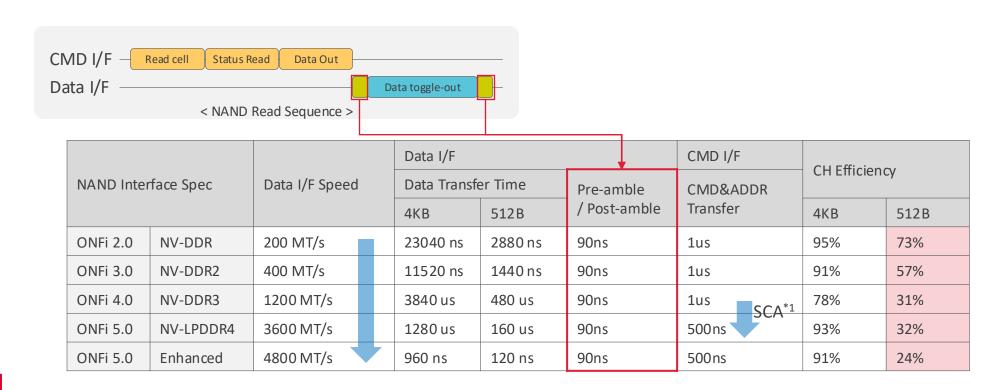


^{*1.} SCA: Separated Command Address

Why the NAND Channel is the Bottleneck

Factors that reduce channel efficiency

- 1. Slow improvement of the CMD I/F speed with legacy manner
- 2. Fixed overhead for data interfaces reliability
 - Pre-amble, Post-amble

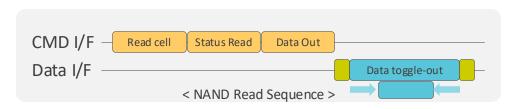


^{*1.} SCA: Separated Command Address

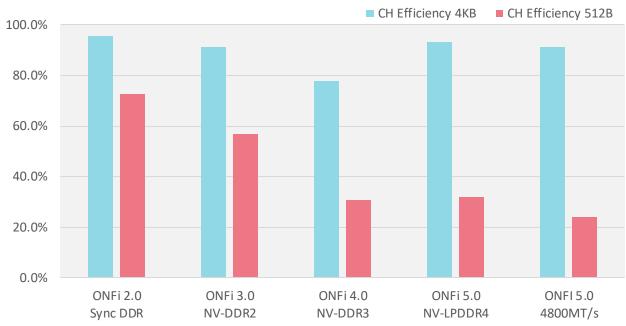
Why the NAND Channel is the Bottleneck

Factors that reduce channel efficiency

- 1. Slow improvement of the CMD I/F speed with legacy manner
- 2. Fixed overhead for data interfaces reliability
 - Pre-amble, Post-amble
- 3. The smaller the data size, the relatively higher the fixed cost



Channel Efficiency with Data I/F Speed & Data Size



NAND Interface Spec D				Data I F Fixed Cost		CMD I/F	CH Efficiency		
		Data I/F Speed Dat		Data Transfe	Data Transfer Time Pre-amble C				CMD&ADDR
				4KB	512B	/ Post-amble	Transfer	4KB	512B
ONFi 2.0	NV-DDR	200 MT/s		23040 ns	2880 ns	90ns	1us	95%	73%
ONFi 3.0	NV-DDR2	400 MT/s		11520 ns	1440 ns	90ns	1us	91%	57%
ONFi 4.0	NV-DDR3	1200 MT/s		3840 us	480 us	90ns	1us SCA*1	78%	31%
ONFi 5.0	NV-LPDDR4	3600 MT/s		1280 us	160 us	90ns	500ns	93%	32%
ONFi 5.0	Enhanced	4800 MT/s		960 ns	120 ns	90ns	500ns	91%	24%

*1. SCA: Separated Command Address

Points to optimize for 100M IOPS

Goal to optimize

- Channel Efficiency of 75% or higher
- 6250K IOPS per Channel (= 160 ns per single read command)
 - CMD I/F occupancy time < 160 ns
 - Data I/F occupancy time < 160 ns

2 Blockers

CMD I/F time: 500ns > 160 ns
 Data I/F time: 210ns > 160 ns

			Data I/F		CMD I/F	CII Efficience		
NAND Inter	face Spec	Data I/F Speed	Data Transfe	er Time	Pre-amble	CMD&ADDR	CH Efficien	-y
			4KB	512B	/ Post-amble	Transfer	4KB	512B
ONFi 5.0	Enhanced	4800 MT/s	960 ns	120 ns	90ns	500ns	91%	24%

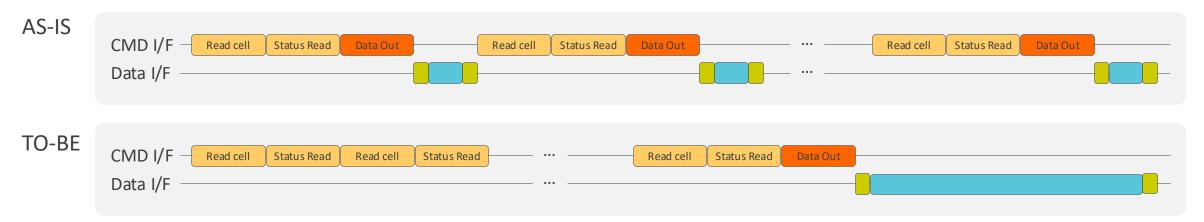
Limitation

Pre-amble & Post-amble time is fixed

Approach to solve problem

Reducing the fixed overhead ratio

Data size	① Data-out Time	② Pre-amble +Post-amble time	Overhead Ratio (= 2 / (1 + 2))	Data I/F efficienty
512B	119ns	90ns	43%	57%
4K	955ns	90ns	8.6%	91.4%
16K	3820ns	90ns	2.3%	97.7%



Limitation

Pre-amble & Post-amble time is fixed

Approach to solve problem

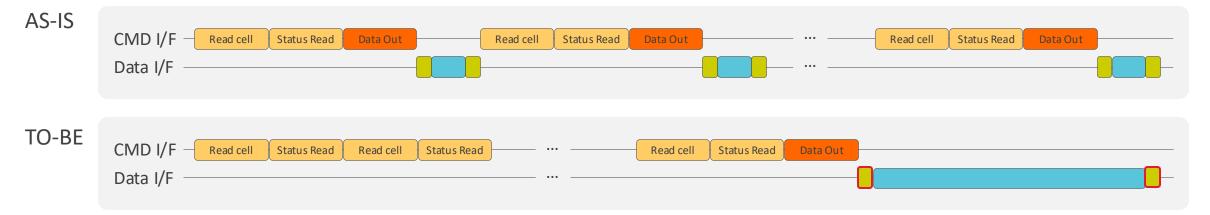
Reducing the fixed overhead ratio

Data size	① Data-out Time	2 Pre-amble +Post-amble time	Overhead Ratio $(= 2)/(1+2)$	Data I/F efficiency
512B	119ns	90ns	43%	57%
4K	955ns	90ns	8.6%	91.4%
16K	3820ns	90ns	2.3%	97.7%

Improvement Effect

Data I/F

- Overhead sharing



Limitation

Pre-amble & Post-amble time is fixed

Approach to solve problem

Reducing the fixed overhead ratio

Data size	① Data-out Time	② Pre-amble +Post-amble time	Overhead Ratio $(= 2)/(1+2)$	Data I/F efficiency
512B	119ns	90ns	43%	57%
4K	955ns	90ns	8.6%	91.4%
16K	3820ns	90ns	2.3%	97.7%

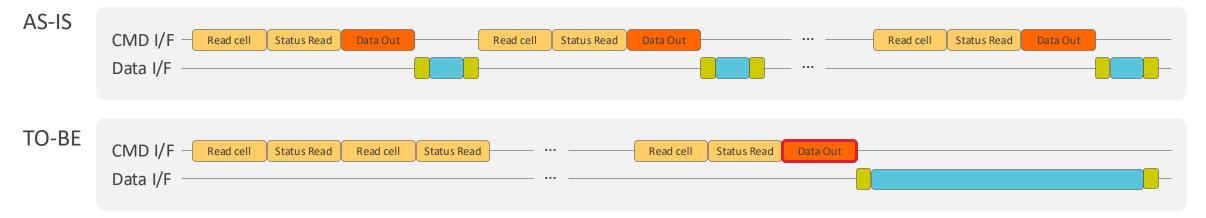
Improvement Effect

Data I/F

Overhead sharing

CMD I/F

 Reduction in the number of Data Out commands



Limitation

Pre-amble & Post-amble time is fixed

Approach to solve problem

Reducing the fixed overhead ratio

Data size	① Data-out Time	2 Pre-amble +Post-amble time	Overhead Ratio $(= 2)/(1+2)$	Data I/F efficiency
512B	119ns	90ns	43%	57%
4K	955ns	90ns	8.6%	91.4%
16K	3820ns	90ns	2.3%	97.7%

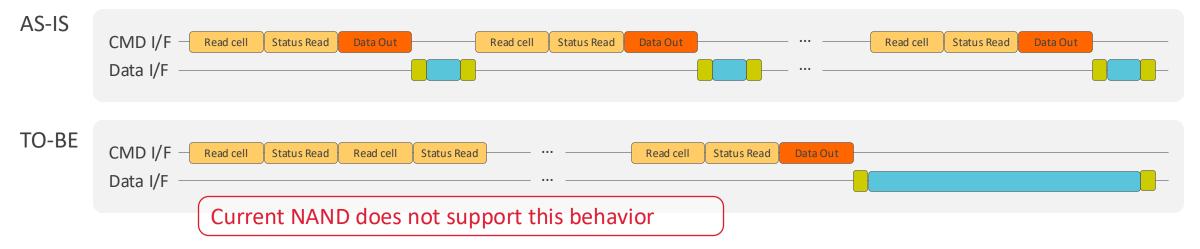
Improvement Effect

Data I/F

- Overhead sharing

CMD I/F

 Reduction in the number of Data Out commands

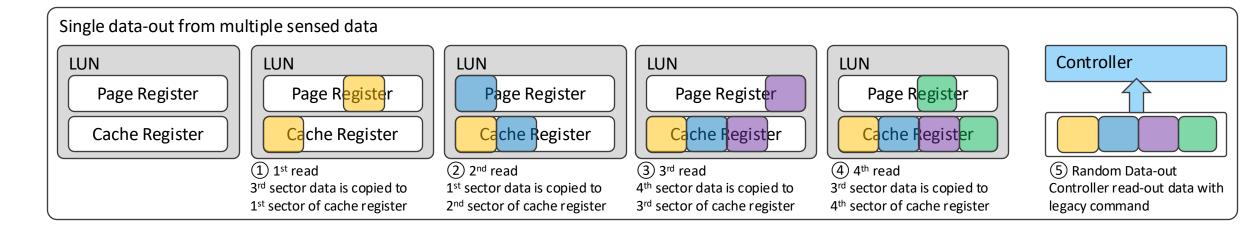




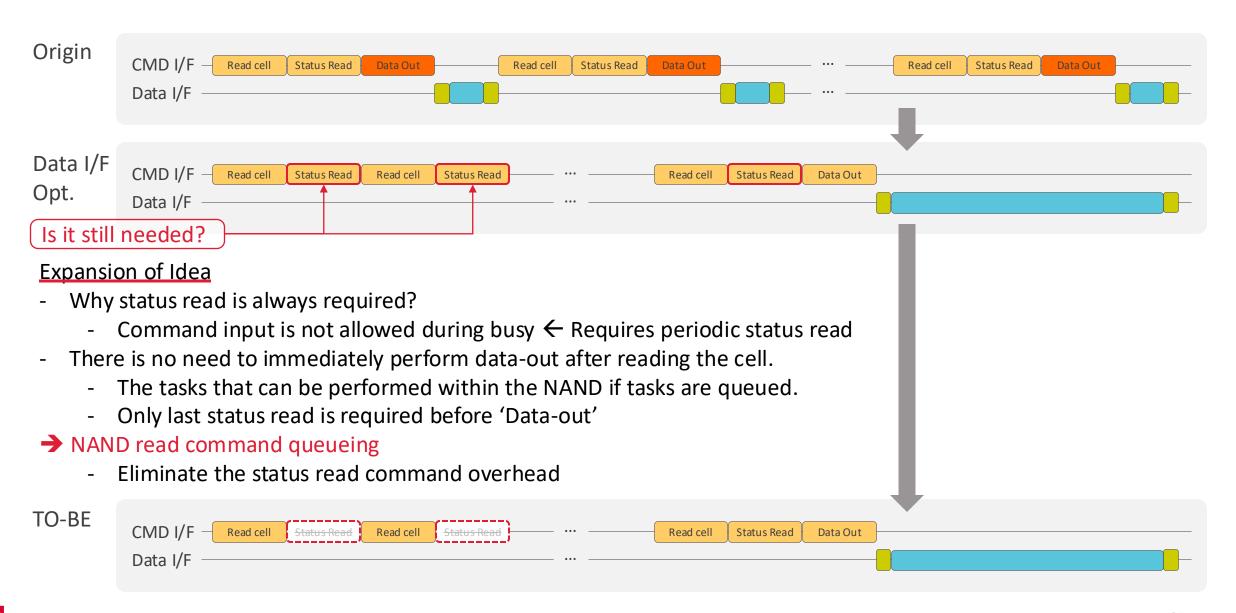
Required New NAND Command To-Be:

Read & Flexible load to cache register

- After cell sensing, data is copied to cache register with source & destination offset
- Cache register is utilized like data FIFO
- Improve the constraint requiring loading to a fixed column address between the page register and cache register



NAND Channel Optimization – CMD I/F 1st

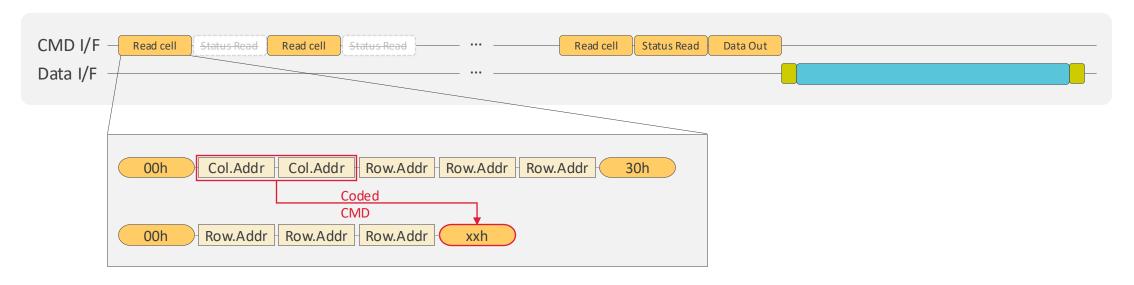


NAND Channel Optimization – CMD I/F 2nd

Idea to reduce CMD I/F overhead

Coded CMD for Limited Column Address Usage

- The usage case for NAND column addresses is aligned to ECC units.
- 2Byte information is not required.
- Legacy confirm command(30h) is replaced by coded command include column address info(ex. 30h~4Fh)
- → Reduction of CMD I/F occupancy time to send read command



Conclusion

- The customized SSD market is rapidly growing.
 - To meet the requirements of the custom SSD market, improvements are needed at each layer.
- SSD solution & NAND vendors must jointly enhance the customized solutions.
 - FADU is engaging to advance the future solutions with multiple NAND suppliers .
- The NAND Legacy CMD I/F is somewhat unsuitable for AI SSD.
 - The fixed overhead ratio will continue to increase.
- Long-term, the NAND Command Interface requires innovative improvements.

+

