

Small Granularity Graph Neural Network Training and the Future of Storage

September 15th, 2025

John Mazzie

MTS, Systems Performance Engineer

Why do we care about storage?

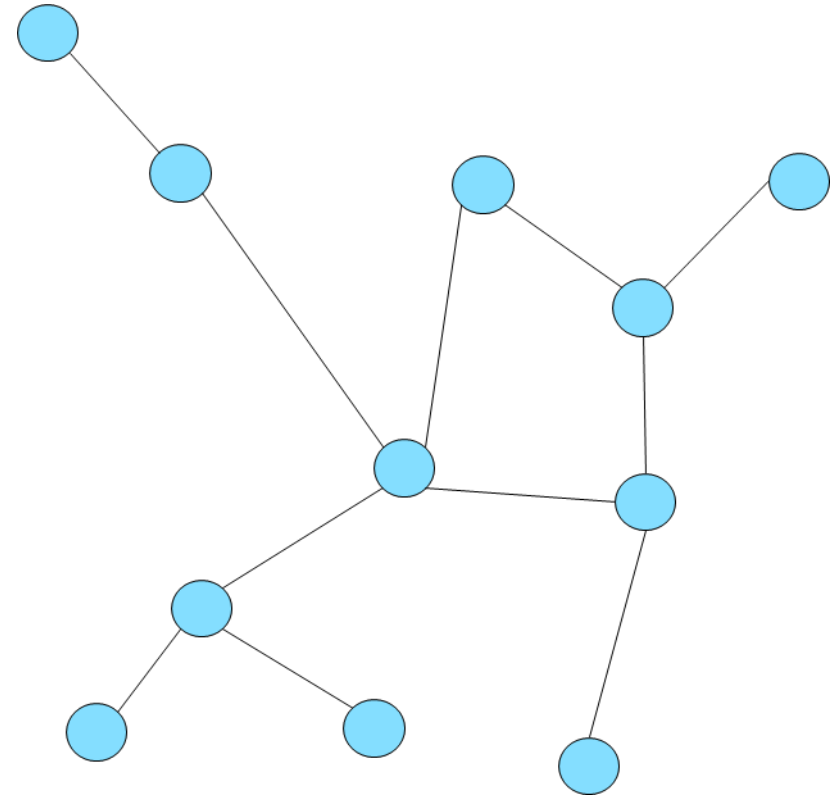
Why do we care about storage?

- Models are getting too large to fit in GPU and System Memory
- Storage must be used in some way for large datasets during workload training/execution
- Model Examples
 - Illinois Graph Benchmark Dataset
 - Graph Neural Network Model
 - Heterogenous 600M nodes
 - 2.3TB on Disk
 - Meta Llama4 Maveric
 - Large Language Inference Model
 - 402B Total Parameters
 - 800GB on Disk

Graph Neural Networks

Graph Structure

- Graph Structure Data
 - Nodes – Entity within graph
 - Edges – Relationship between nodes
 - Stored in sparse matrix formats
- Node Feature Data
 - Feature embeddings for each node
 - Stored in $N \times D$ matrix
 - N – Number of nodes
 - D – Dimension of feature



GNN Training Process

- Mini-batching
 - Splitting graphs into smaller subgraphs
- Graph sampling
 - Choosing which nodes to sample
- Feature Aggregation
 - Gathering of node features of subgraph
- Data Transfer
 - Transfer of features to GPU memory
- Model Training
 - Magic!

Feature Aggregation Using GIDS

Large Datasets

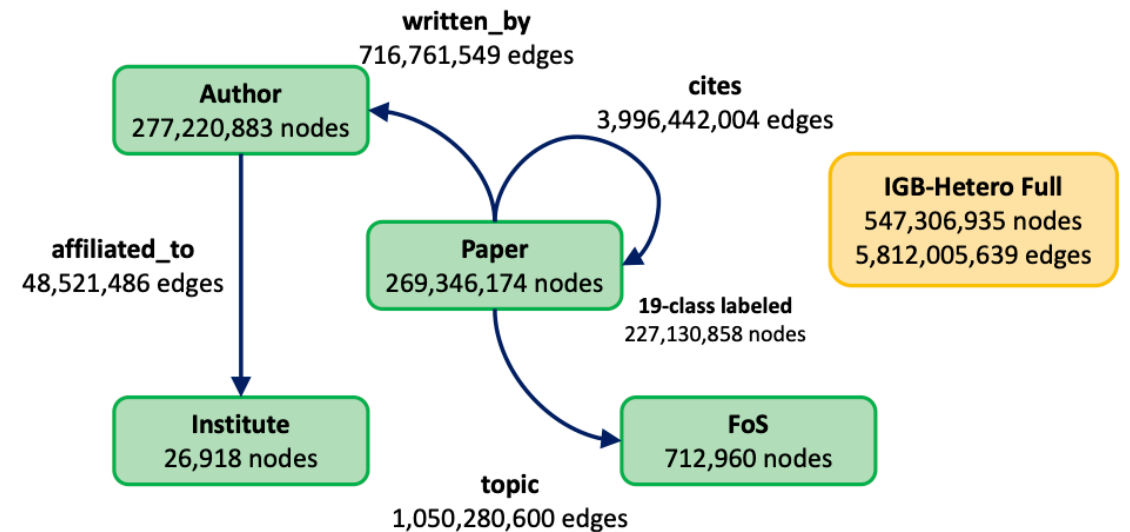
- Current State of the art
 - Memory map storage to system memory
 - High CPU overhead
 - Low Performance
- GIDS (GPU Initiated Direct Storage)
 - Prototype project from NVIDIA Research
 - Proprietary NVMe driver
 - GPU coordinates all I/O transfers
 - CPU is completely bypassed (for feature aggregation)
 - Leverages high parallelism of GPU
 - High Performance

Dataset and Performance

Illinois Graph Benchmark (IGB) Dataset

Evaluation Graph

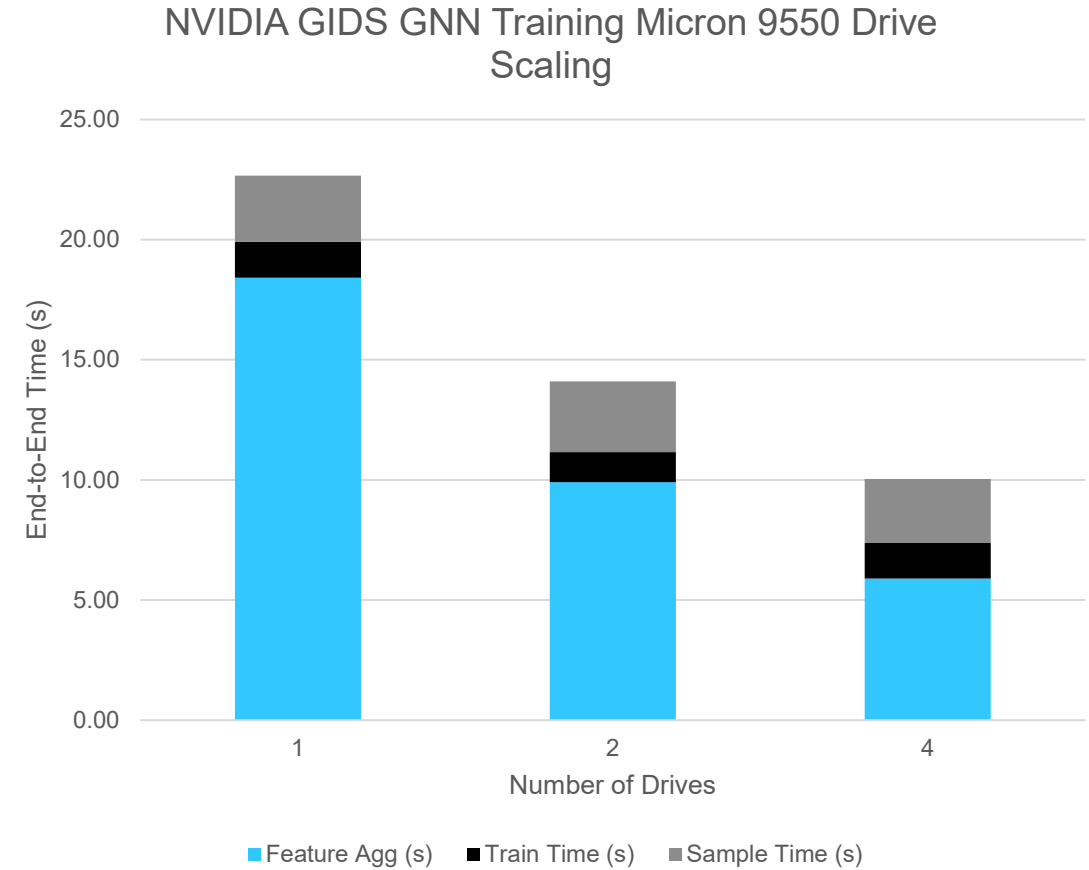
- Motivation
 - Lack of large-scale datasets
 - Lack of labeled data
- IGB
 - Homogenous and heterogenous options
 - Large number of labeled nodes
 - Variable sized embeddings
 - Variable sized graphs
 - 4KiB features



GIDS GNN Training Results

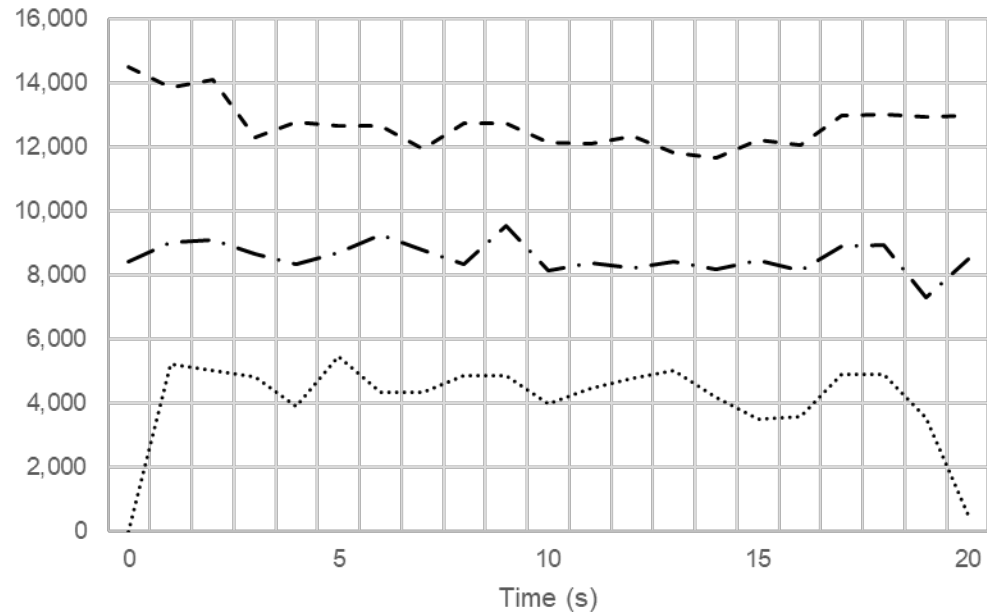
4096 Batch Size - 8GiB Cache Size

	4x 9550s	2x 9550s	1x 9550	MMAP (Gen4 NVme)
Sampling	2.67	2.94	2.76	4.65
Feature Aggregation	5.89	9.9	18.41	1,130.12
Training	1.48	1.26	1.49	2.13
End-to-End	10.04	14.1	22.61	1,142.90





GIDS GNN Training Queue Depth over Time

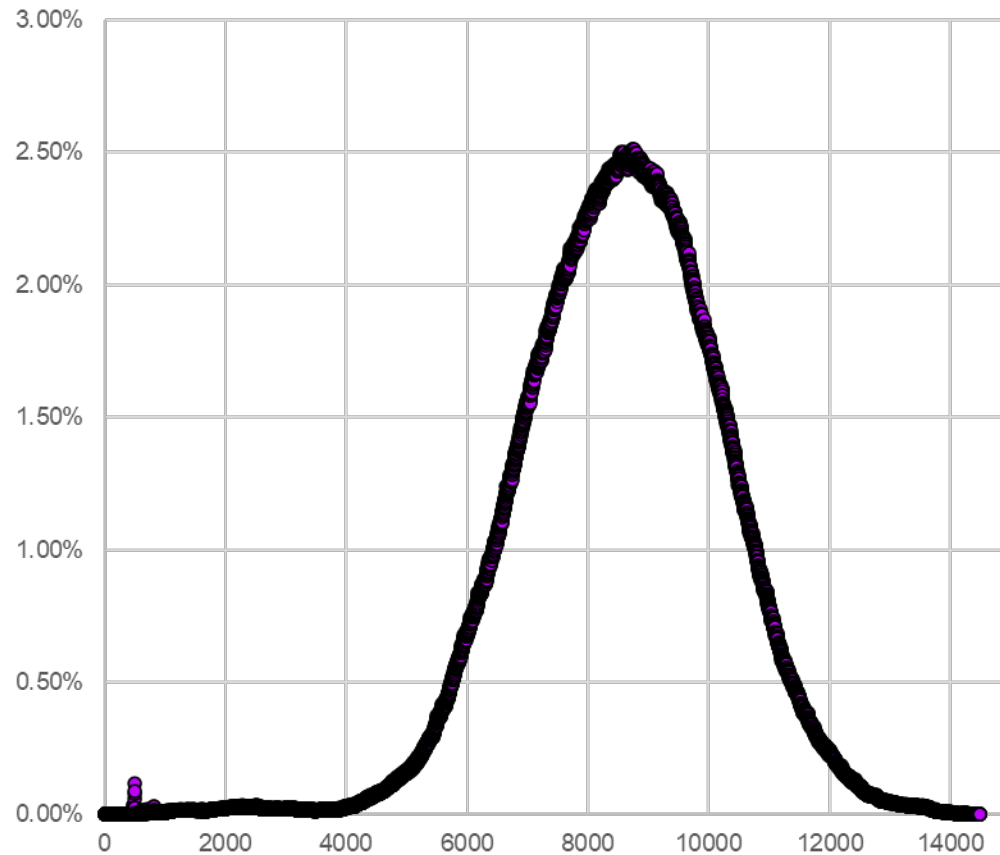


- Typical Enterprise Application
 - Order of 100 QD is considered high
- GIDS GNN Training
 - Minimum: Between 4K and 6K
 - Average: Between 8K and 10K
 - Peak: Between 12K and 14K

..... Minimum Queue Depth - . - Average Queue Depth - - - Maximum Queue Depth



GIDS GNN Training Queue Depth Histogram



- Most of the time queue depth is between 6K and 12K
- Much higher queue depth than expected when designing NVMe devices
- Will this workload benefit from more queue pairs?
- What other optimizations can we make to handle these higher number of transactions?

SSD considerations for GNN use case

Small Transaction Handling of SSDs

- **Full Page Read Internally:**
 - The SSD reads the entire page (Typically 4KiB) from NAND into its internal DRAM or SRAM buffer.
- **Data Extraction:**
 - The SSD controller extracts only the requested portion (e.g., 512B) from the buffer and sends it to the host.
- **Caching and Optimization:**
 - SSDs often use **read-ahead** and **caching** strategies to anticipate future reads, improving performance.
 - If multiple small reads target the same page, the SSD can serve them from cache without re-reading NAND.
- **Impact on Performance:**
 - While sub-4KB reads are supported, they can be **less efficient** due to:
 - Overhead of reading full pages.
 - Increased wear if small writes follow (due to read-modify-write cycles).
 - High IOPS workloads with small reads can still perform well, especially on enterprise SSDs with optimized firmware.

Synthetic Testing of Smaller Transactions

- Hardware
 - PCIe Gen5 NVMe SSD x4
 - PCIe Gen5 and Gen4 based servers
- Software
 - nvm-block-bench
 - Synthetic benchmark included in Big Accelerator Memory (BaM)

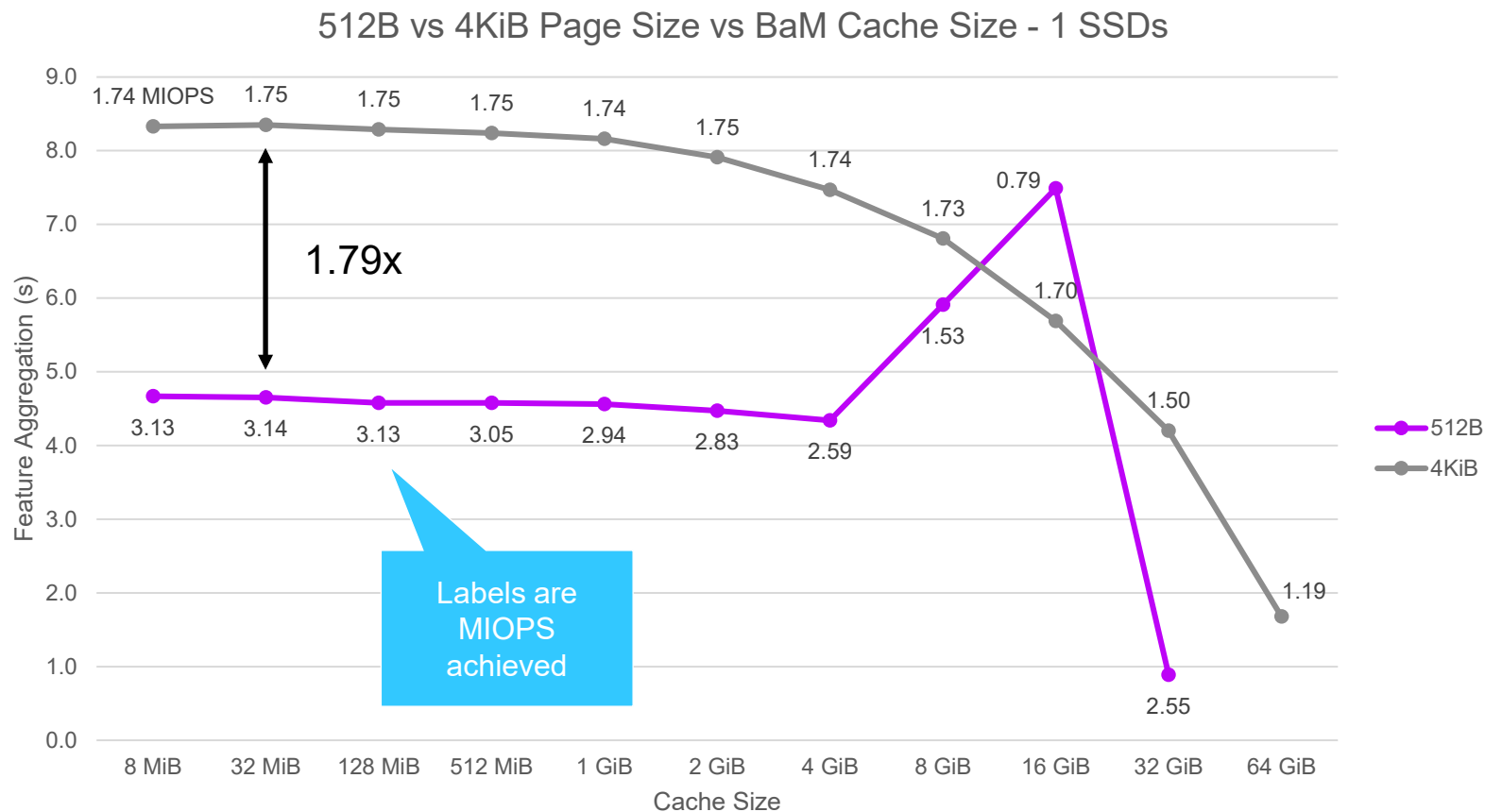
PCIe Gen5 SSD				
	PCIe Gen4 Bus		PCIe Gen5 Bus	
Page Size (B)	IOPS	GiB/s	IOPS	GiB/s
4096	6.9M	26.32	13.5M	51.81
2048	13.5M	25.76	13.5M	25.87
1024	13.5M	12.94	13.5M	12.95
512	13.5M	6.46	13.5M	6.3

Small Transaction GNN Training

- Training on modified version heterogenous graph used in previous testing
 - Features modified from 4KiB to 512B
- System Configuration
 - Test system is a PCIe Gen4 System with NVIDIA A100
 - Drives are Micron 9550 NVMe devices (PCIe Gen5)
- What are we comparing
 - GNN Training w/ BaM cache using 4KiB pages
 - GNN Training w/ BaM cache using 512B pages
 - BaM cache sizes

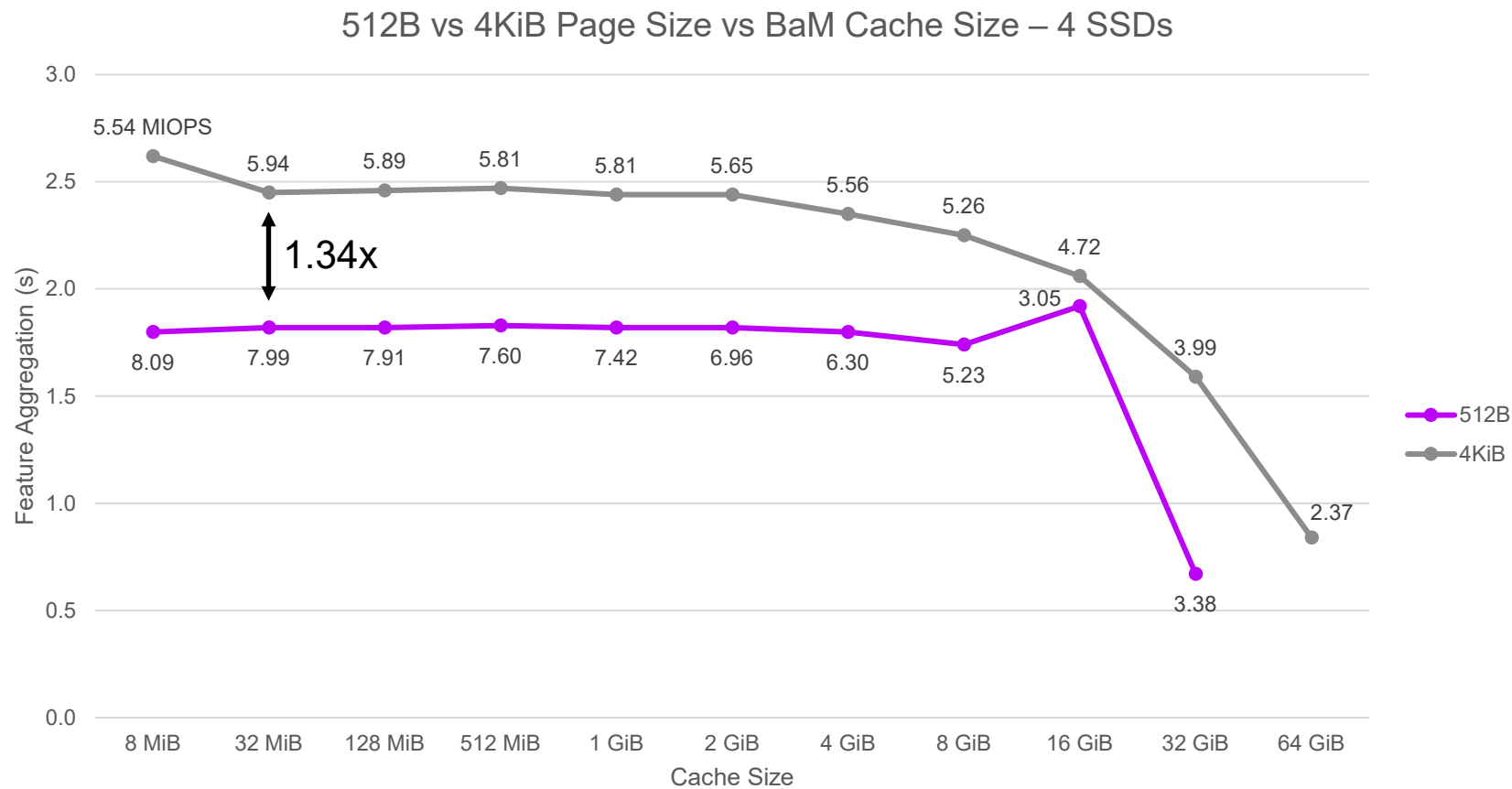
Small Transaction GNN Training – 1 NVMe

- Feature Aggregation improvement with 512B transactions
- Takes advantage of Gen5 device IOPS, while on a Gen4 interface
- Shows potential of improvement with new devices



Small Transaction GNN Training – 4 NVMe

- Don't get 2x improvement like the single drive case.
- 4-8x more 512B transactions on a device would be ideal.



Ecosystem Enablement

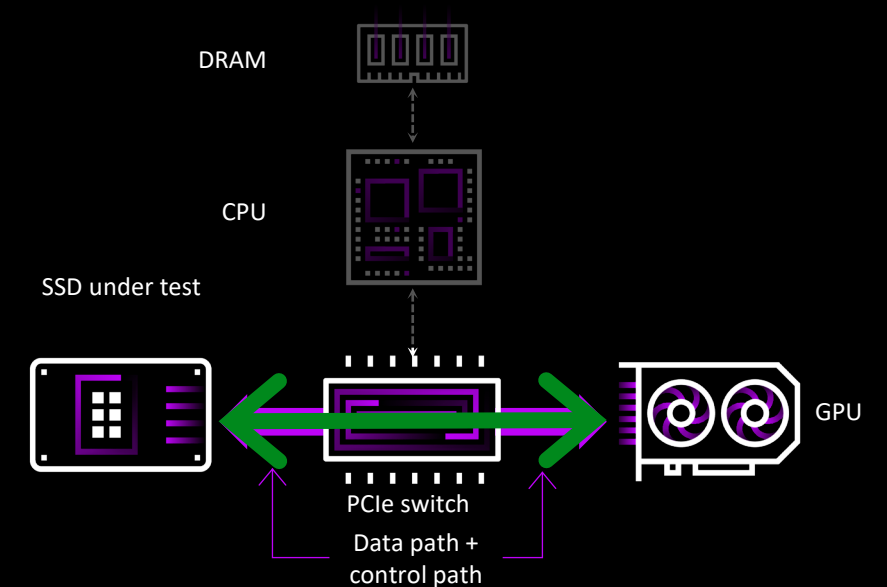
Shift in Storage Paradigm

A GPU can drive high levels of IO traffic

- Max IOPs on 1 CPU core is ~1M
- 100M IOPs = 100 cores, just for IO
- AI accelerators have tens of thousands of cores and can use them for massively parallel IO
- Faster SSDs are needed to keep up with AI workload demands

GPU Initiated Storage

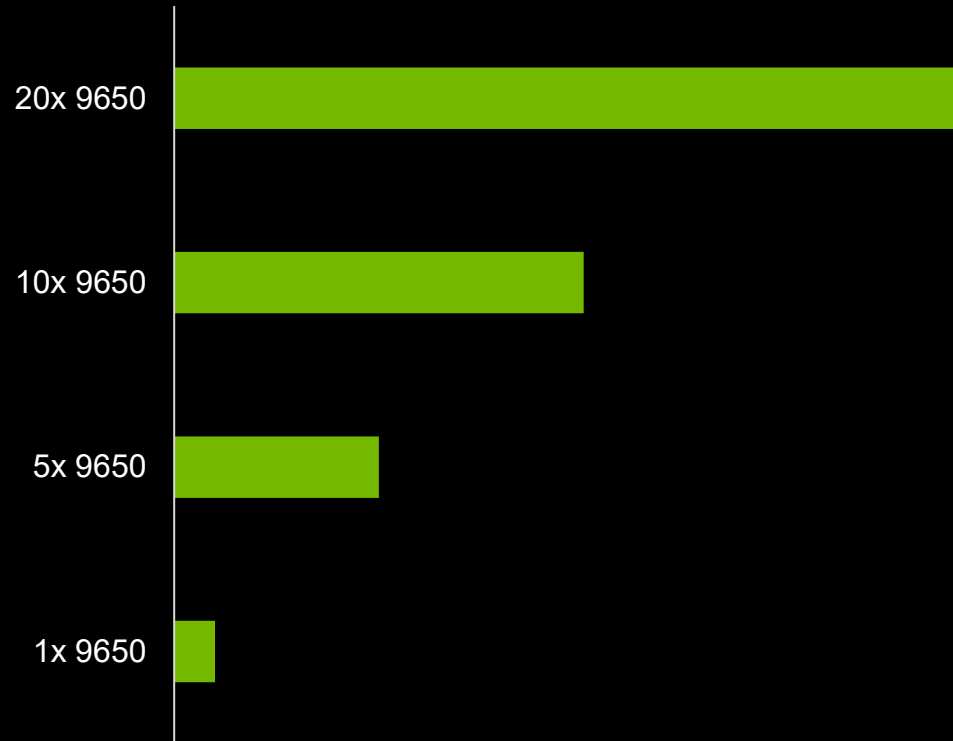
- **NVIDIA GDS:** Data flows direct to GPU, control plane travels over CPU+DRAM
- **Big Accelerator Memory (BAM):** Data and control plane traffic flow directly to GPU over PCIe switch complex
- **NVIDIA SCADA:** Like BAM, plus client server architecture and advanced features



Micron 9650 | NVIDIA SCADA

Early SCADA code shows strong performance and linear scaling on H3 Platform (Preliminary Results)

NVIDIA SCADA (Preview Software)



Early NVIDIA SCADA code drives impressive small block random read IOPs through 20 Micron 9650 Gen6 NVMe SSDs

- **Linear performance scaling from 1 to 20 drives**
- H3 Platform System:
 - Intel 8568Y+, 512GB DDR5
 - 3x Broadcom 144 lane PCIe Gen6 Switches
 - 20x Micron 9650 Gen6 NVMe SSD, E1.S 7.68TB
 - H100 NVL 96GB HBM3
- * Preliminary Results:
 - * Hardware & software stack tuning for ongoing

SCADA test results collected on pre-production code.
System under test: H3 Platform PoC system, 1x Intel 8568Y+, 512GB DDR5, 3x Broadcom A0 PCIe Gen6 switches, 20x Micron 9650 E1.S 7.68TB, NVIDIA H100NVL-96GB PCIe Gen5x16,
Workload is 512B random read initiated from H100 GPU.
Performance testing completed by Micron's Data Center Workload Engineering team.

Storage Next

- Organized by NVIDIA with NDA stakeholder partners
- Forward-looking initiative aimed at redefining storage architectures for AI workloads
- Industry participants from NAND and controller vendors, storage providers, CSPs, and OEMs
- Focus
 - GPU Initiated I/O
 - High IOPS
 - Small Data (Sub-4KiB)
 - Power Efficiency



© 2025 Micron Technology, Inc. All rights reserved. Information, products, and/or specifications are subject to change without notice. All information is provided on an "AS IS" basis without warranties of any kind. Statements regarding products, including statements regarding product features, availability, functionality, or compatibility, are provided for informational purposes only and do not modify the warranty, if any, applicable to any product. Drawings may not be to scale. Micron, the Micron logo, the M logo, Intelligence Accelerated™, and other Micron trademarks are the property of Micron Technology, Inc. All other trademarks are the property of their respective owners.