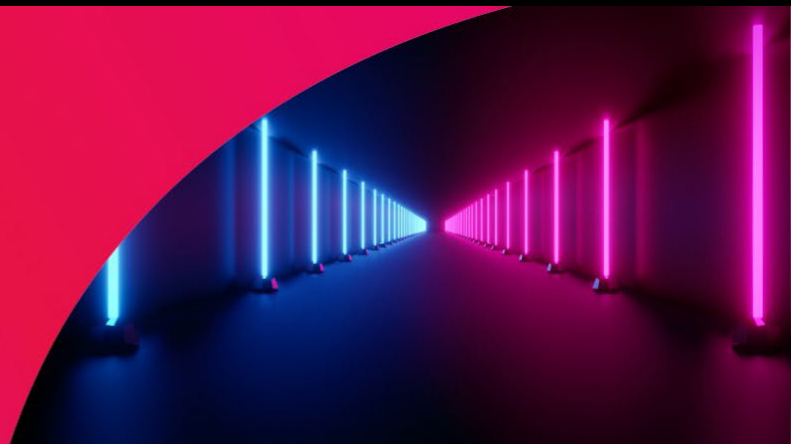


Micron® 9550 NVMe™ SSD performance with NVIDIA® Magnum IO™ GPUDirect® Storage



Imagine a superhighway built for speed and efficiency—a direct route that connects GPU data demands directly to NVMe™ SSD storage.

That's what NVIDIA® GPUDirect® Storage (GDS) is in the world of AI. Part of the powerful framework of NVIDIA's Magnum IO™ acceleration technologies, GDS is a game-changer supported with the Micron 9550 NVMe SSD.¹

The result? A dramatic surge in system bandwidth combined with a significant reduction in power used across different training workloads.²

This document analyzes:

- Measured throughput (performance in GB/s)
- SSD power efficiency (in GB/s per watt)
- Energy consumed (in joules) to transfer a fixed amount of data (energy efficiency)

Testing was performed across three transfer sizes: 4KB, 128KB and 1MB for the Micron 9550 SSD, the Kioxia CM7-R, and the Samsung PM1743 (all SSDs' advertised capacity = 7.68TB).³

Test results show the Micron 9550 SSD enables greater storage performance for AI workloads and is more performant and energy-efficient compared to these NVMe SSDs.



Figure 1: Micron 9550 SSD (U.2; 800GB to 30.72TB)

Key findings

As AI models have grown rapidly in size and complexity, managing their power consumption has become increasingly difficult.

Testing showed the Micron 9550 SSD enabled faster data transfer at 4KB, 128KB and 1MB transfer sizes, demonstrated up to 34% higher throughput, up to 70% better power efficiency, and consumed 81% less energy when transferring a 1 TB data set.

34% Faster throughput

The Micron 9550 SSD showed faster data throughput (GB/s) at 4KB, 128KB and 1MB transfer sizes compared to the competitors' SSD. The maximum improvement of 34% was measured using a 4KB transfer size.

76% Better power efficiency

Head-to-head testing demonstrated the Micron 9550 SSD offered up to 76% better power efficiency (measured in GB/s per watt). Its maximum advantage was observed with a 128KB transfer size.

81% Less energy used

Energy consumption is a critical factor to consider when building AI systems. One way to help reduce energy consumption is to build more energy-efficient hardware, like the Micron 9550 SSD. When transferring 1TB of data, the Micron 9550 SSD used up to 81% less energy than the competitors' SSDs.

micron.com/9550

1. See the [NVIDIA GPUDirect Storage Overview Guide](#) for additional information on GDS.
2. Micron internal engineering analysis of AI training workloads shows that different IO sizes are seen depending on model and data formats. Therefore, this document focuses on small (4KB), medium (128KB), and large (1MB) transfer sizes in two test scenarios.
3. Unformatted. 1 GB = 1 billion bytes. Formatted capacity is less. Competitive PCIe Gen5 SSDs (Kioxia CM7-R and Samsung PM1743) chosen from the top 10 PCIe SSD suppliers shown in the Forward Insights analyst report "SSD Supplier Status Q1/24 May 2024".

GPUDirect® Storage – a direct path between GPUs and NVMe SSDs⁴

GPUDirect Storage (GDS) creates a direct data path between local or remote storage, such as NVMe (or NVMe over Fabrics [NVMe-oF]; this discussion focuses on local NVMe) and GPU memory. By enabling a direct-memory access (DMA) engine near the network adapter or storage, data is moved into or out of GPU memory—without burdening the CPU (the control path still relies on the CPU). This direct IO path is represented in Figure 2.

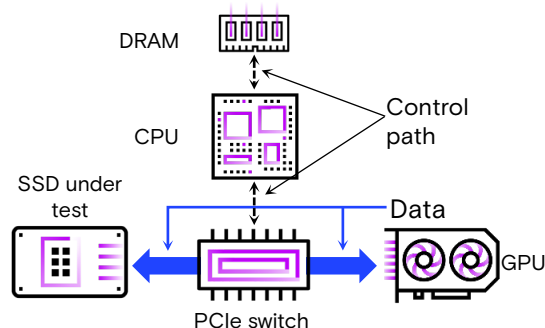


Figure 2: GDS IO path

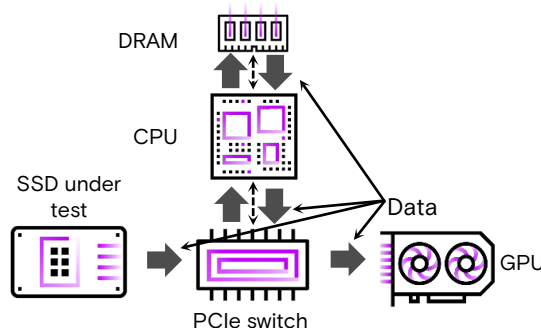


Figure 3: Legacy IO path

This more direct data path between storage and GPU memory is designed to avoid extra copies through a bounce buffer in the CPU's memory required by a legacy IO path represented in Figure 3. All tests use the GDSIO tool to measure results using the GDS IO path.⁵

4KB transfers: Up to 34% higher performance and up to 76% better power efficiency

With over 16,000 cores, the NVIDIA H100 Tensor Core GPU drives high parallelism in AI workloads. In testing GDS on an AI system, the number of GDSIO workers was increased until reaching maximum SSD performance. Figure 4 represents performance results (in GB/s) with a 4KB workload. The Micron 9550 SSD is shown in purple while the Kioxia CM7-R is shown in grey and the Samsung PM1743 is shown in dark grey.

Each SSD shows maximum performance at 512 GDSIO workers.⁶ The Micron 9550 SSD shows 30% higher performance than the Kioxia CM7-R and 34% higher performance than the Samsung PM1743. Figure 5 represents power efficiency in GB/s per watt, with efficiency data at 512 GDSIO workers is highlighted. Here, the Micron 9550 shows 73% higher power efficiency than the Kioxia CM7-R, and 76% higher power efficiency than the Samsung PM1743.⁷

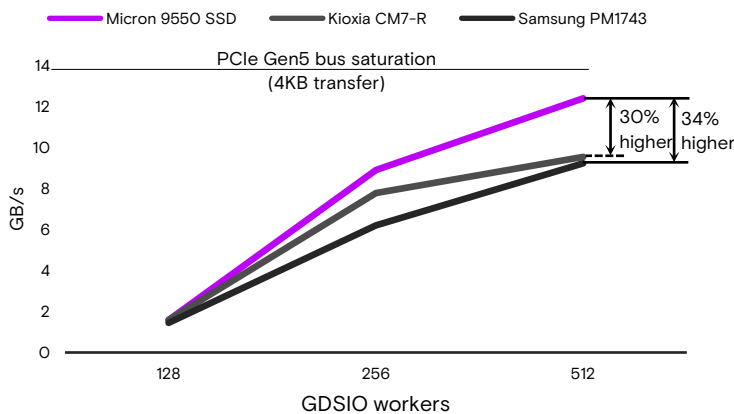


Figure 4: 4KB transfer performance (GB/s, higher is better)

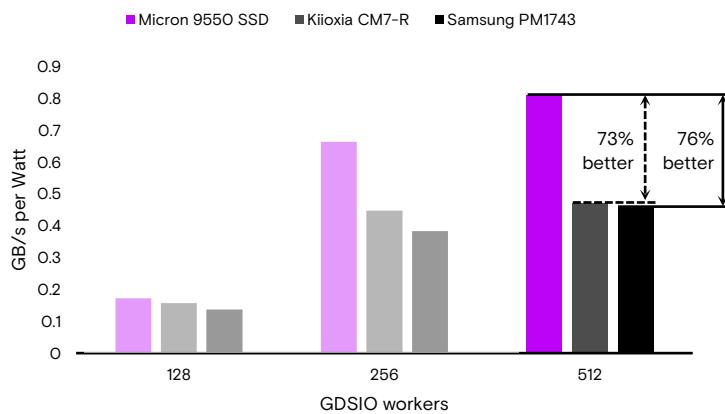


Figure 5: 4KB transfer power efficiency (GB/s per watt, higher is better)

4. See <https://developer.nvidia.com/gpudirect-storage> for additional details on the IO path differences.

5. GDSIO is a tool that simulates increasing IO parallelism and bandwidth requirements of a GPU. See [NVIDIA GPUDirect Storage Benchmarking and Configuration Guide](#) for additional information on this tool.

6. Black solid line indicates PCIe Gen5 x4 bus saturation of 13.92GB/s for 4KB transfer size (assuming MPS of 512B and common link and common host efficiencies).

7. Performance and power efficiency differences are calculated as (higher value / lower value) - 1, expressed as a percentage. Power data for 512 GDSIO workers is highlighted in Figure 5. In each of the subsequent power efficiency figures, the highlighted value is the number of GDSIO workers at highest performance.

128KB transfers: Up to 25% higher performance and up to 70% better power efficiency

Figure 6 represents performance results with a 128KB workload. The Micron 9550 SSD is shown in purple while Kioxia CM7-R is shown in grey and the Samsung PM1743 is shown in dark grey. Each SSD shows maximum performance at 64 GDSIO workers. The Micron 9550 SSD shows 21% higher performance than the Kioxia CM7-R and 25% higher performance than the Samsung PM1743.⁸

Figure 7 represents power efficiency (in GB/s per watt) with power efficiency data for 64 GDSIO workers highlighted. The Micron 9550 shows 70% better power efficiency than the Kioxia CM7-R and 61% better power efficiency than the Samsung PM1743.

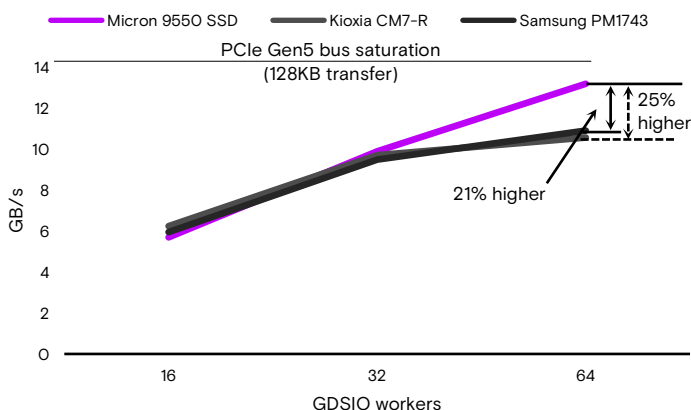


Figure 6: 128KB transfer performance (GB/s, higher is better)

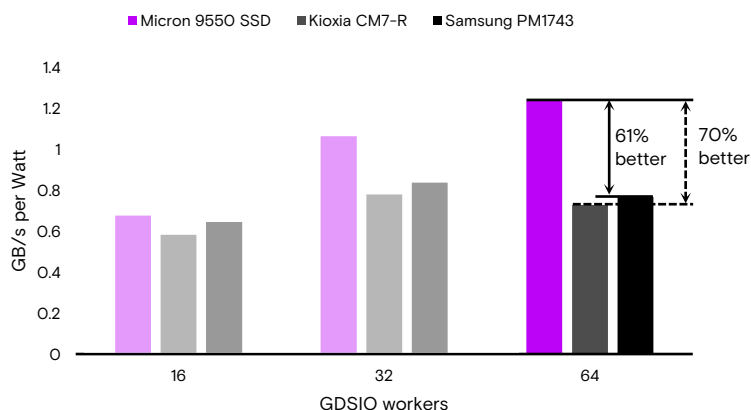


Figure 7: 128KB transfer power efficiency (GB/s per watt, higher is better)

1MB transfers: Up to 14% higher performance and up to 19% better power efficiency

Figure 8 represents performance results in GB/s with a 1MB workload. The Micron 9550 SSD is shown in purple while the Kioxia CM7-R is shown in grey and the Samsung PM1743 is shown in dark grey. Each SSD shows maximum performance at 16 GDSIO workers. The Micron 9550 SSD shows 14% higher performance than the Kioxia CM7-R and 9% higher performance than the Samsung PM1743.⁹

Figure 9 represents power efficiency (in GB/s per watt) with power efficiency data for 16 GDSIO workers highlighted. The Micron 9550 shows 19% better power efficiency than the Kioxia CM7-R and 12% better power efficiency than the Samsung PM1743.

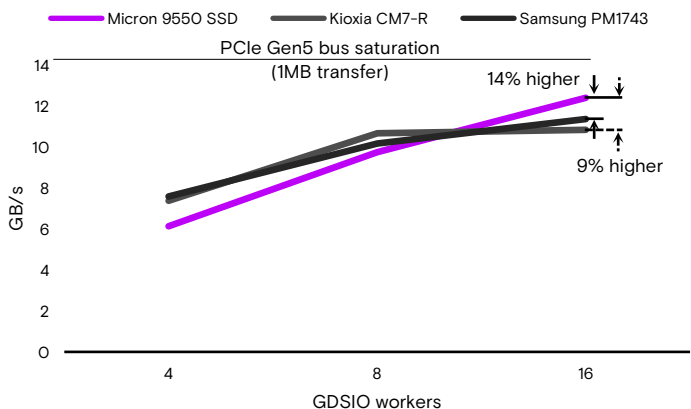


Figure 8: 1MB transfer performance (higher is better)

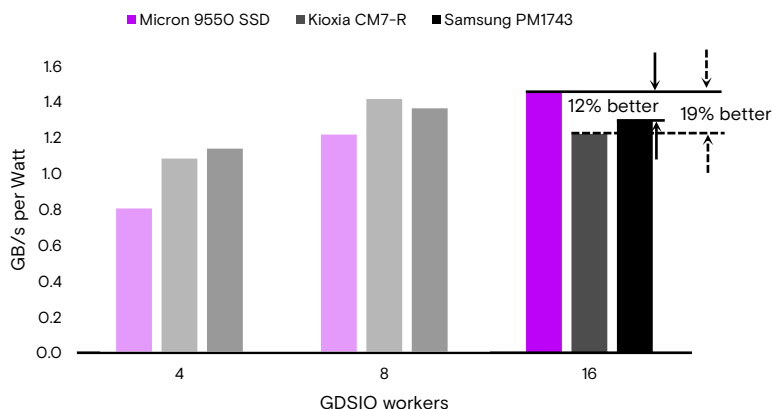


Figure 9: 1MB transfer power efficiency (higher is better)

8. Black solid line indicates PCIe Gen5 x4 bus saturation of 14.31 GB/s for 128KB transfer size (assuming MPS of 512B and common link and host efficiencies).

9. Black solid line indicates PCIe Gen5 x4 bus saturation of 14.32 GB/s for 1MB transfer size (assuming MPS of 512B and common link and host efficiencies).

How to improve energy efficiency

The energy required for AI tasks is growing at an annual rate between 26% and 36%, and by 2028, AI might consume more power than Iceland did in 2021.¹⁰

Because of the associated environmental concerns and operating cost, building power-efficient AI systems is crucial for sustainable AI.

While AI systems are like any other data center systems in many ways, AI energy consumption and its growth put additional emphasis on energy efficiency.

Micron 9550 SSD consumes less energy for every transfer size

Improving efficient hardware can help reduce the environmental impact and operational costs of AI.

SSD energy efficiency is one such factor. It can be compared by examining the relative energy consumption for each SSD (at each transfer size's maximum performance) for a fixed amount of data transferred.

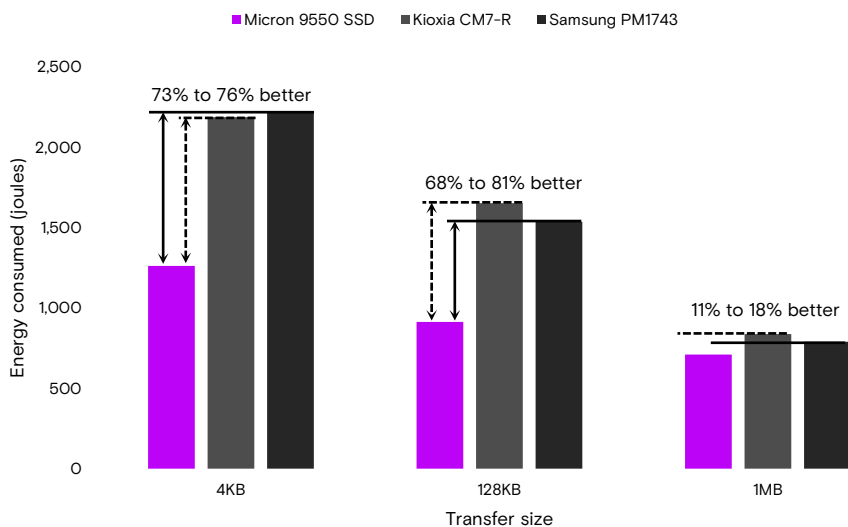


Figure 10: Average SSD power consumption for each transfer size (lower is better)

Figure 10 represents the energy consumed (in joules, 1 joule = 1 watt-second) along the vertical axis for each SSD at each transfer size (along the horizontal axis). The values in Figure 10 use maximum SSD performance noted earlier to transfer 1TB of data. Table 1 summarizes the results.^{11,12}

Transfer size	Micron 9550 SSD energy (joules)	Kioxia CM7-R energy (joules)	Micron 9550 SSD advantage	Samsung PM1743 energy (joules)	Micron 9550 SSD advantage
4KB	1,260	2,185	73%	2,214	76%
128KB	913	1,650	81%	1,536	68%
1MB	709	837	18%	790	11%

Table 1: Relative SSD energy used to transfer 1TB of data by transfer size

10. See [How to manage AI's energy demand – today, tomorrow and in the future](#) for additional details on relative energy consumption.

11. Energy consumption advantage percentage calculated as (comp. joules) – (Micron 9550 SSD joules) / (comp. joules), expressed as a percentage.

12. Energy is defined as power multiplied by time; thus, both SSD performance and power use are factors.

Conclusion

Using the Micron 9550 NVMe SSD with GDS shows that the SSD offers significant performance and efficiency advantages compared to the Kioxia CM7-R and Samsung PM1743 SSDs. The Micron 9550 SSD demonstrated up to 34% higher performance and up to 76% better power efficiency compared to these SSDs as well.

In terms of energy efficiency, the Micron 9550 SSD showed up to 81% better energy efficiency when transferring large data sets.

How we tested

Tables 2 and 3 outline the Supermicro with NVIDIA H100 system and software configurations used.

Supermicro server with NVIDIA H100 configuration	
Server	Supermicro® SYS-521GE-TNRT
CPU	2X Intel® Xeon® Platinum 8568Y+ 48-core processors
Memory	16X Micron 96GB DDR4 RDIMMs @ 5600MT/s (1.5TB total memory)
Network	1X NVIDIA® H100-NVL 96GB
Micron SSDs	Micron 7.68TB 9550 NVMe SSD
Competitors' SSDs	Kioxia CM7-R, Samsung PM1743 (7.68TB advertised capacity for each SSD)

Table 2: Hardware configuration

Software configuration	
OS	Ubuntu 20.04.6 LTS
Kernel	5.4.0-177-generic
CUDA	12.4
NVIDIA Mellanox OFED version	24.01-0.3.3
NVIDIA GDSIO tool	1.11
Filesystem	XFS
Data Layout	4GB file per GDSIO worker and 1,024 total files equates to 4TB used

Table 3: Software configuration

©2024 Micron Technology, Inc. All rights reserved. All information herein is provided on as “AS IS” basis without warranties of any kind, including any implied warranties, warranties of merchantability or warranties of fitness for a particular purpose. Micron, the Micron logo, and all other Micron trademarks are the property of Micron Technology, Inc. All other trademarks are the property of their respective owners. Products are warranted only to meet Micron’s production data sheet specifications. Products, programs and specifications are subject to change without notice. Rev. A 07/2024 CCM004-676576390-11761