

Maximizing exascale AI workloads with the Micron 6550 ION NVMe SSD



At present, data centers worldwide consume 1-2% of overall power, but according to Goldman Sachs Research, this percentage will likely rise to 3-4% by the end of the decade.¹

Power-efficient SSDs like the Micron® 6550 ION SSD reduce data center power consumption to enable allocating more power to GPUs, CPUs, and networking without compromising performance.

This document compares the performance and power improvement of the Micron 6550 ION SSD relative to the Solidigm D5-P5336 in three AI use cases. The Micron 6550 ION SSD's power was limited to 20 watts (PS1), while the Solidigm D5-P5336 could use up to 25 watts (PS0).²

AI checkpointing: Saving the state of a model at regular intervals during training is called checkpointing. This allows the training to resume from the last saved state in case of interruptions or failures.³

DLIO Unet3D: These results help to evaluate how well an SSD performs when handling large 3D medical images for AI training. They also help in understanding SSD performance and power efficiency, which helps ensure the SSD can meet the demanding needs of AI workloads.⁴

DLRM preprocessing: This mechanism transforms raw data into a format suitable for training recommendation models. Test results show how different SSDs help with data size reduction and data set optimizing processes for efficient training.⁵



Figure 1: Micron 6550 ION SSD in U.2 (15mm), E3.S 1T (7.5mm), and E1.L (9.5mm)

Efficient, performant storage for AI

The Micron 6500 ION SSD redefines power-efficient storage for AI: it maximizes IT budgets amidst data growth, performance expectations, and environmental concerns. It is the world's fastest, most energy-efficient 60TB PCIe Gen5 SSD.⁶

AI checkpointing: Compared to the Micron 6550 ION SSD, the Solidigm D5-P5336 takes up to 151% more time, consumes up to 209% more SSD energy, and is up to 135% slower during checkpointing.

- **Less time:** The Micron 6550 ION SSD takes less time to capture AI workload checkpoints.
- **Less energy:** It also takes less energy to capture every checkpoint
- **Better performance:** And it offers superior SSD performance

DLIO Unet3D: The Micron 6550 SSD demonstrates up to 20% better power efficiency and 30% higher performance versus the competitor.

Better power efficiency

Higher performance

20%

30%

DLRM pre-processing: The Micron 6550 SSD takes less energy to complete this workload as well.

The Micron 6550 ION SSD offers significant advantages in both performance and power efficiency. As AI workload power consumption continues to climb, the Micron 6550 ION SSD is an ideal choice to help reduce overall power consumption in data centers while maintaining superior performance levels.

micron.com/6550ION

1. See [this page on goldmansachs.com](https://www.goldmansachs.com) for additional details on this demand increase.
 2. NVMe power states limit the maximum power an SSD can consume (actual power consumption depends on multiple factors). To learn more about NVMe power states, see [this page on nvmexpress.org](https://www.nvmexpress.org).
 3. For additional information on checkpointing, see [this article on restack.io](https://www.restack.io).
 4. For additional information about DLIO and Unet3D, see [this page on the dlio-benchmark website](https://www.dlio-benchmark.com).
 5. See [this page on github.com](https://github.com) for additional information on DLRM pre-processing.
 6. Comparisons are made against the Solidigm D5-P5336. These comparisons use publicly available competitor information from published sources at the time of the 6550 ION launch, with the 6550 ION using a maximum power of 20W and the Solidigm D5-P5336 using a maximum of 25W, resulting in up to 20% less maximum power consumption for the 6550 ION. Improvements calculated as (SSD1 metric / SSD 2 metric) - 1, expressed as a percentage.

AI checkpointing analysis

AI model training checkpointing is the process of saving the state of a model at regular intervals during training, allowing the training to resume from the last saved state in case of interruptions or failures. This technique is crucial for reducing wasted computational resources and ensuring the continuity and efficiency of training large models.⁷

Checkpointing brings several benefits including:

- **Saves points to help with fault recovery:** By checkpointing, system failure risks can be minimized. With checkpoints, extended training time (days or even months) is not lost due to node failure. The training job can be restarted from the last checkpoint after the node is repaired and online.
- **Simple training pause and reuse:** Checkpoints allow pause and resume activity by an AI application at any step during training.
- **Model reuse:** By saving intermediate checkpoints, models can be reused, repurposed, or serve as seeds for new training objectives.

Checkpointing results

Figures 2 through 4 represent the time required for one checkpoint, the SSD energy consumed for one checkpoint, and the SSD throughput during checkpoint operations.⁸

The Micron 6550 ION SSD in PS1 (20 watts maximum) is shown in purple while Solidigm D5-P5336 is in PS0 (25 watts maximum) and shown in dark gray. Setting the Micron 6550 ION SSD to PS1 limits its maximum power draw to 20% less than the Solidigm D5-P5336 in PS0.

Compared to the Micron 6550 ION SSD, the Solidigm D5-P5336 takes up to 151% more time, consumes up to 209% more SSD energy, and is up to 135% slower during checkpointing.

■ Micron 6550 ION SSD (PS1) ■ Solidigm D5-P5336 (PS0)

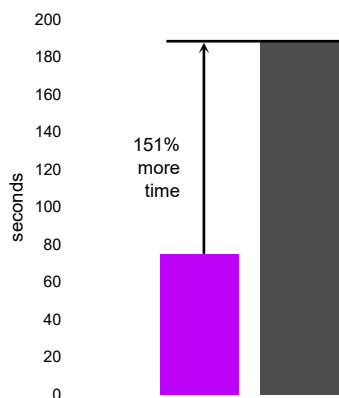


Figure 2: Checkpoint completion time (Micron 6550 ION SSD at 20% lower maximum power)

■ Micron 6550 ION SSD (PS1) ■ Solidigm D5-P5336 (PS0)

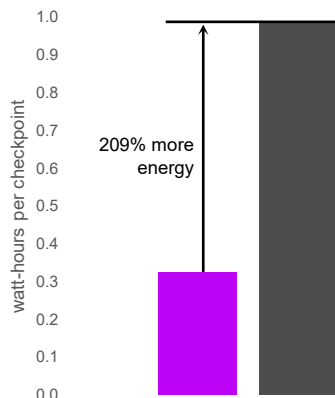


Figure 3: SSD energy to complete checkpoint (Micron 6550 ION SSD at 20% lower maximum power)

■ Micron 6550 ION SSD (PS1) ■ Solidigm D5-P5336 (PS0)

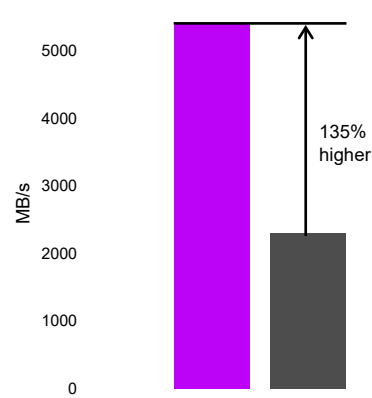


Figure 4: SSD performance in MB/s (Micron 6550 ION SSD at 20% lower maximum power)

The Micron 6550 ION SSD shows faster checkpoint completion time (in seconds), uses less energy to complete a checkpoint (in watt-hours), and higher performance (in MB/s) than the Solidigm D5-P5336, despite an intentional 20% lower maximum power limitation on the Micron SSD.

In Figure 2, checkpoint completion time in seconds is shown on the vertical axis (lower is better). Figure 2 shows that the competitor SSD takes 151% more time to complete a checkpoint than the Micron 6550 ION SSD, and Figure 3 shows that the competitor SSD consumes 209% more energy per checkpoint. Figure 4 shows SSD performance (in MB/s) along the vertical axis (higher is better) during checkpointing and that the competitor SSD is up to 135% slower than the Micron 6550 ION SSD.

7. See [this page on deepchecks.com](https://www.deepchecks.com) for more details on checkpointing ---and its value.

8. Checkpoint workload modeled on Llama3 405B parameter LLM. Model representing an 8GPU server. Checkpoint size is 415GB.

DLIO Unet3D analysis

Deep learning I/O (DLIO) Unet3D benchmark is a method used to evaluate SSD performance in medical imaging using a model that reads large image files into accelerator memory and generates dense volumetric segmentations. It employs the same data loaders as those in real workloads, such as PyTorch and TensorFlow, to move data from storage to CPU memory.⁹

This analysis helps in understanding storage system performance characteristics when handling complex tasks. It provides insights into the storage system’s throughput, power efficiency, and overall performance. This information is crucial for optimizing storage solutions for AI training, ensuring that the storage can manage the demanding data requirements of AI workloads.¹⁰

Higher performance and power efficiency are beneficial as the former helps the system process data more quickly and the latter can help reduce a data center’s environmental impact.

Unet3D results

Figures 5 and 6 represent SSD performance and power efficiency results using the Unet3D benchmark with three simulated H100 accelerators.

In Figures 5 and 6, the Micron 6550 ION SSD in PS1 (20 watt maximum) is shown in purple while Solidigm D5-P5336 in Power State 0, PS0, (25 watt maximum) is shown in dark grey. Setting the Micron 6550 ION SSD to PS1 limits its maximum power draw to 20% lower than the Solidigm D5-P5336 in PS0.

The Micron 6550 ION SSD delivers up to 20% better power efficiency and 30% SSD higher performance

The Micron 6550 ION SSD shows better SSD power efficiency (in MB/s per watt) and higher performance (in MB/s) than the Solidigm D5-P5336, despite an intentional 20% lower maximum power limitation on the Micron SSD.

In Figure 5, power efficiency in MB/s per watt is shown on the vertical axis (higher is better). Figure 5 shows that the Micron 6550 ION SSD delivered 20% better power efficiency, even when its maximum power consumption was limited to just 20 watts.

Figure 6 shows SSD performance (in MB/s) along the vertical axis (higher is better) under the same maximum power conditions noted in Figure 5. The Micron 6550 ION SSD delivered 30% higher performance — even when its maximum power consumption is restricted to 20% lower than the Solidigm D5-P5336 maximum.

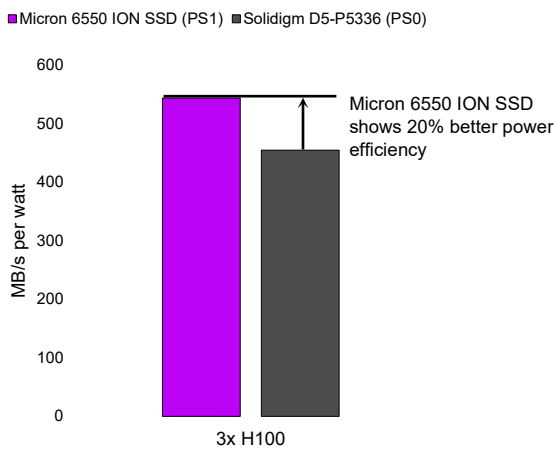


Figure 5: SSD power efficiency (Micron 6550 ION SSD at 20% lower maximum power)

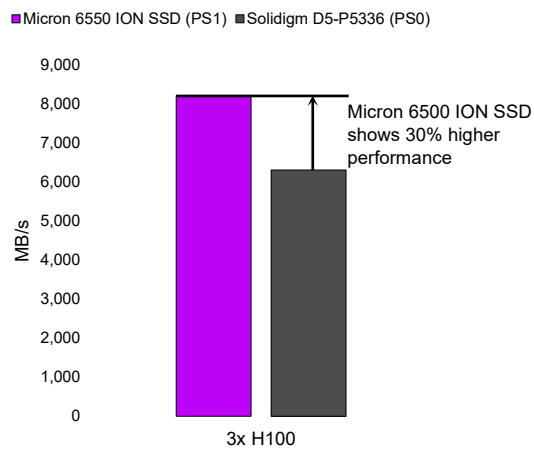


Figure 6: SSD performance (Micron 6550 ION SSD at 20% lower maximum power)

9. See [the DLIO page on github](#) for more information on this benchmark and its use.

10. See [this page on mlcommons.org](#) for additional information on the benefits of high-performance, power-efficient SSDs in generative AI storage.

DLRM and pre-processing analysis

Deep Learning Recommendation Model (DLRM) pre-processing is a crucial step in preparing data for training recommendation models. This process involves converting raw data into a format suitable for model training, typically by transforming the data into a columnar format, like Parquet, and then into a binary format.

The pre-processing stage is essential because it ensures that the data is clean, well-structured, and optimized for efficient training. For instance, in the context of DLRM, pre-processing can significantly reduce the data size from several terabytes to a more manageable size, making it easier to handle and process. The DLRM pre-processing completion time and energy efficiency are both important to AI use cases. Superior energy efficiency (measured as SSD watt-seconds necessary to complete the workload) helps consume less energy to complete the workload. Because this workload is highly CPU- and DRAM-dependent, SSD performance is expected to be similar.¹¹

Test results scale from 1 to 24 workers (the number of parallel processes or threads that are employed to handle the pre-processing tasks). Tests conclude when the pre-processing is completed. Because this workload has a fixed endpoint, we can analyze SSD energy used (power in watts multiplied by workload duration in hours to show energy used in watt-hours) rather than just SSD power draw.

Pre-processing results

Because performance (completion time) is expected to be similar between the tested SSDs, energy consumption is the most relevant, storage-centric metric to compare the Micron 6550 SSD (PS1) to the Solidigm D5-P5336 (PS0).

Figure 4 shows SSD energy consumed (in watt-hours) along the vertical axis (lower is better) and the number of workers on the horizontal axis. The Micron 6550 ION SSD is shown in purple, while the Solidigm D5-P5336 SSD is shown in dark gray.

The Solidigm D5-P5336 SSD consumed up to 65% more energy than the Micron 6550 ION SSD (the greatest difference is seen at 24 workers).

Since DLRM pre-processing is heavily CPU- and DRAM-dependent, similar workload completion times are expected for both SSDs.

This can be seen in Figure 5 where workload completion time (in seconds) is shown along the vertical axis and the number of workers is shown along the horizontal axis.

As expected, this workload shows similar completion time for each SSD configuration at each number of workers. Although the competitor SSD shows longer completion times for each worker count compared to the Micron 6550 ION SSD, the differences are minimal.

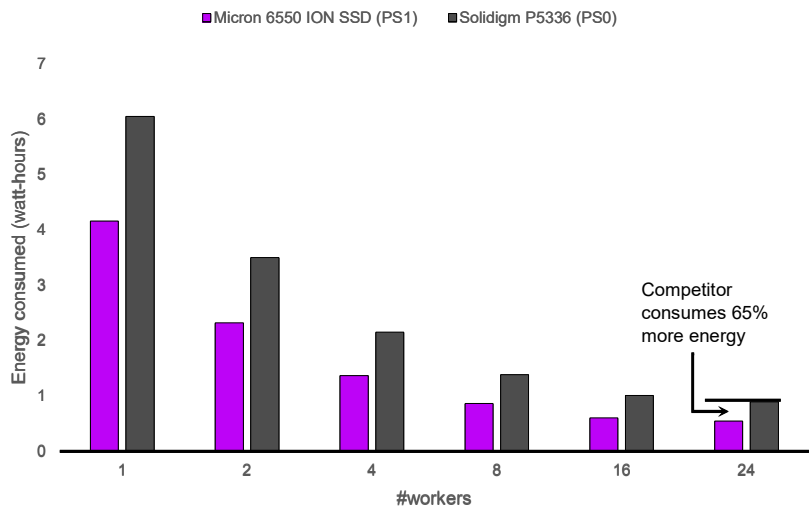


Figure 4: DLRM pre-processing workload SSD energy consumed (lower is better)

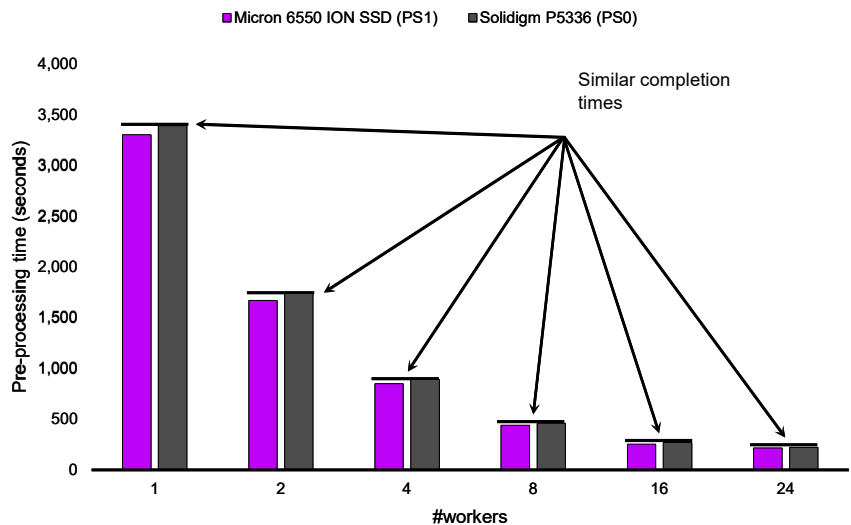


Figure 5: Pre-processing workload completion time in seconds (lower is better)

11. See [this technical paper on arxiv.org](#) and [this github.com page](#) for more information on DLRM dependency on DRAM and CPU resources.

Conclusion

The Micron 6550 ION SSD demonstrates significant advantages in both performance and energy consumed compared to the Solidigm D5-P5336, making the Micron 6550 ION SSD an ideal choice for AI workloads where SSD capacity is paramount.

The Micron 6550 ION SSD demonstrates better capture time and less energy use when checkpointing AI workloads, better energy efficiency and higher performance in DLIO Unet3D testing, as well as better energy consumption results for DLRM pre-processing — all with a maximum power draw that is limited to 20% less than the competitor.

These advantages highlight the Micron 6550 ION SSD's ability to reduce overall energy consumption in data centers, and specifically in AI workloads, allowing more power allocation to critical components like GPUs and CPUs without compromising performance. Overall, the Micron 6550 ION SSD maximizes IT budgets amidst data growth, performance expectations, and environmental concerns.

Visit the [Data center SSD storage page on micron.com](#) to start leveraging the benefits of the Micron 6550 ION SSD.

Server configurations used

	Checkpointing and DLIO Unet3D platform	DLRM pre-processing platform
Server platform	Supernode AS-1115CS-TNR	Supernode AS-1115CS-TNR
CPU	1x AMD EPYC™ 9654	1x AMD EPYC™ 9654
Memory	256GB	768GB
Server storage	1x Micron 6550 ION SSD (61.44TB) 1x Solidigm D5-P5336 (61.44TB)	1x Micron 6550 ION SSD (61.44TB) 1x Solidigm D5-P5336 (61.44TB)
Boot drive	Micron 7450 SSD M.2 (480GB)	Micron 7450 SSD M.2 (480GB)
Operating system	Ubuntu 20.04 LTS (Focal Fossa)	Ubuntu 20.04 LTS (Focal Fossa)
Additional software	Kernel – 5.15.0-105-generic	Kernel – 5.15.0-105-generic Apache PyArrow – 17.0.0 Numpy – 2.0.2 Pandas – 2.0.3

Table 1: Server configurations