

# Analyzing the Effects of GPUDirect Storage on AI Workloads

**SDC EMEA 2021**

Or Lapid, Field Applications Engineer

© 2021 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 orbit logo, the M orbit logo, Intelligence Accelerated™, and other Micron trademarks are the property of Micron Technology, Inc. All other trademarks are the property of their respective owners.



# Agenda

- Can Storage Impact the AI training performance?
- NVIDIA DALI as a proxy for benchmarking maximum throughput
- NVIDIA GPUDirect Storage (GDS)
  - What is it?
  - Micron Tests results of Local NVMe Performance GDS vs. Legacy data-path.



Founded more than 40 years ago on **October 5, 1978**

Headquartered in **Boise, Idaho, USA**

**\$21.4B**

FY2020 annual revenue

**4th**

largest semiconductor company in the world

**134**

on the 2020 Fortune 500

**44,000**

patents granted and growing

**17**

countries

**13**

manufacturing sites and 14 customer labs

**40,000**

team members

# Recap of previous findings

## Can storage impact training performance?

- 8x NVIDIA V100
- Container resource limits are used to show impacts of constrained systems.
- Training speed of ResNet-50 model with ImageNet dataset.

### Memory:

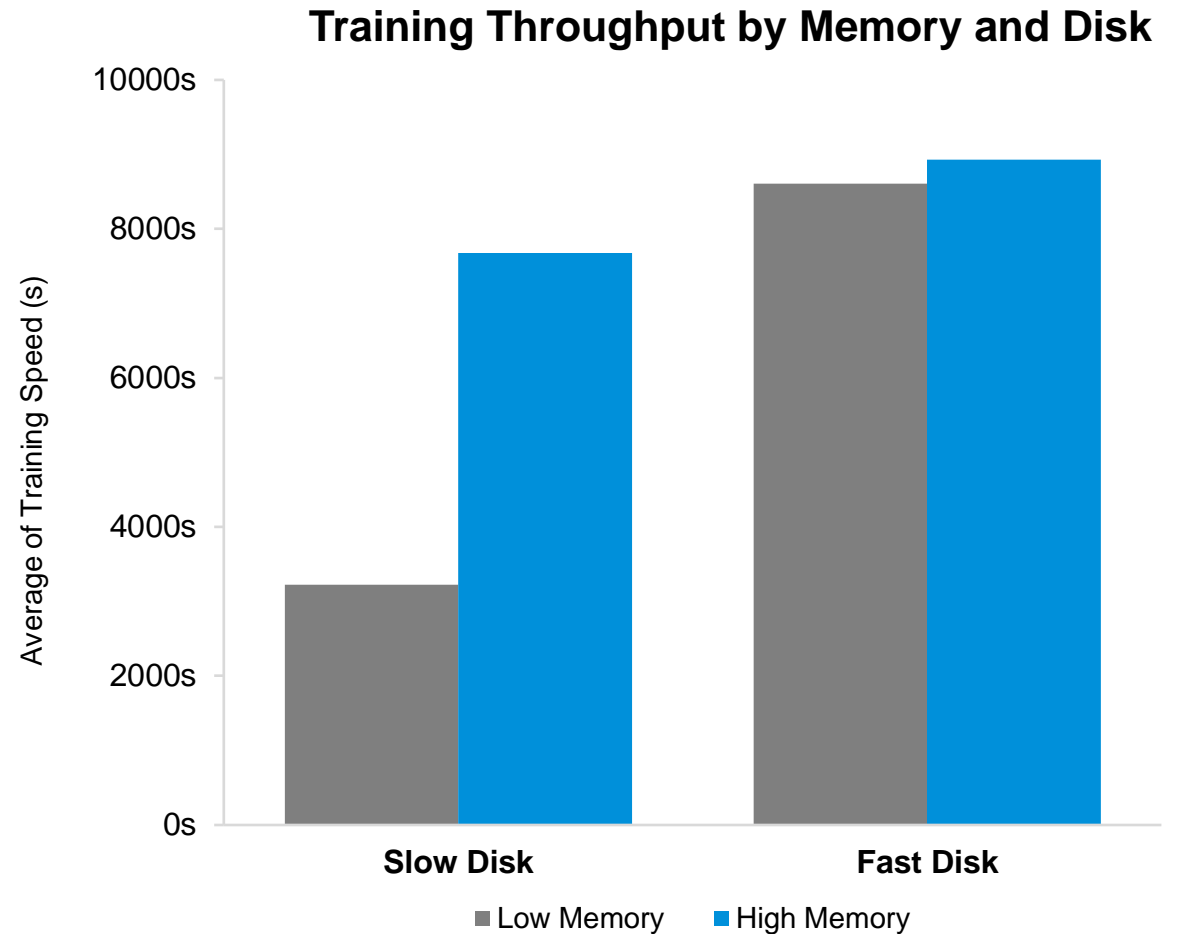
High memory = 1TB

Low Memory = 128GB

### Disk:

Fast Disk = 8x NVMe

Slow Disk = 500MB/s Limit



# NVIDIA Data Loading Library (DALI)

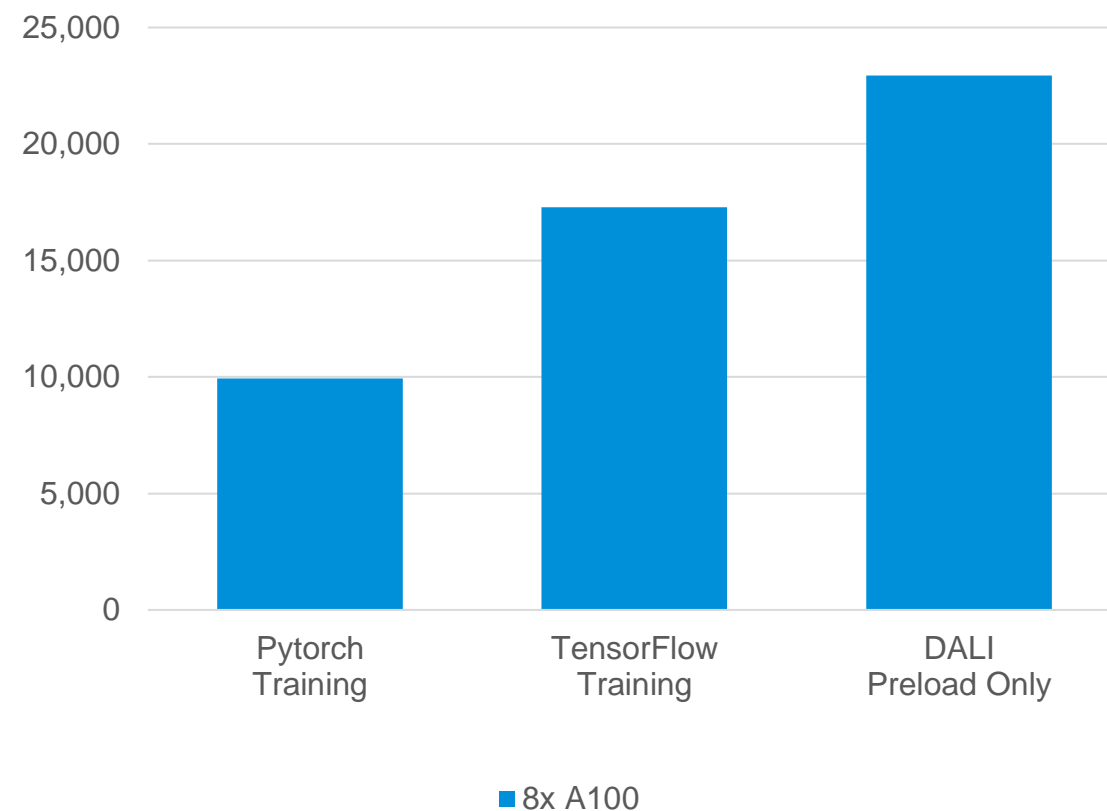
Used here to benchmark the “maximum theoretical” image throughput

- The defined pipeline has 3 steps: Read, Image Decode, Resize to model input
- Container was memory limited to ensure accesses went to storage
- 4 threads per GPU reading from RecordIO file
- Images are the Imagenet 2012 dataset
- Hardware is the NVIDIA DGX A100

Data compared to training throughput with 2 popular frameworks

- ResNet-50 model trained on Imagenet 2012
- Training results from NVIDIA: <https://developer.nvidia.com/deep-learning-performance-training-inference>

Images per Second for Training and Data Preloading



# NVIDIA Data Loading Library (DALI)

Compute intensive training is currently at 75% of “maximum”

Faster storage won't increase the top end performance without architectural changes

Moving data through the CPU memory as a bounce buffer can result in a “storage” bottleneck that can't be fixed with faster storage

Leads us to NVIDIA **GPUDirect Storage**

The data from this benchmark—like most benchmark data—is somewhat contrived.

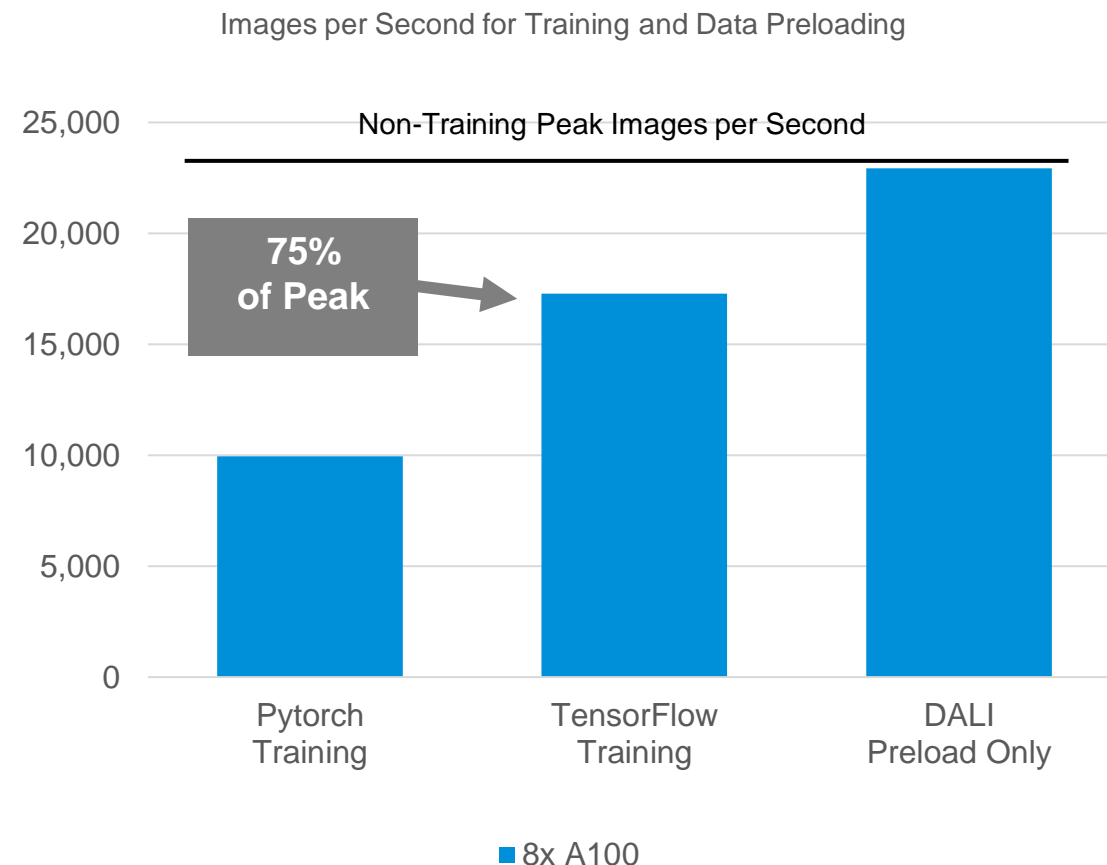
Data layout could be optimized better

DALI preload performance could be improved to about double what we see here.

The data layout and pipeline settings were selected to match NVIDIA's submissions to MLPerf for image classification instead of optimizing purely for performance.

But that sort of misses the point.

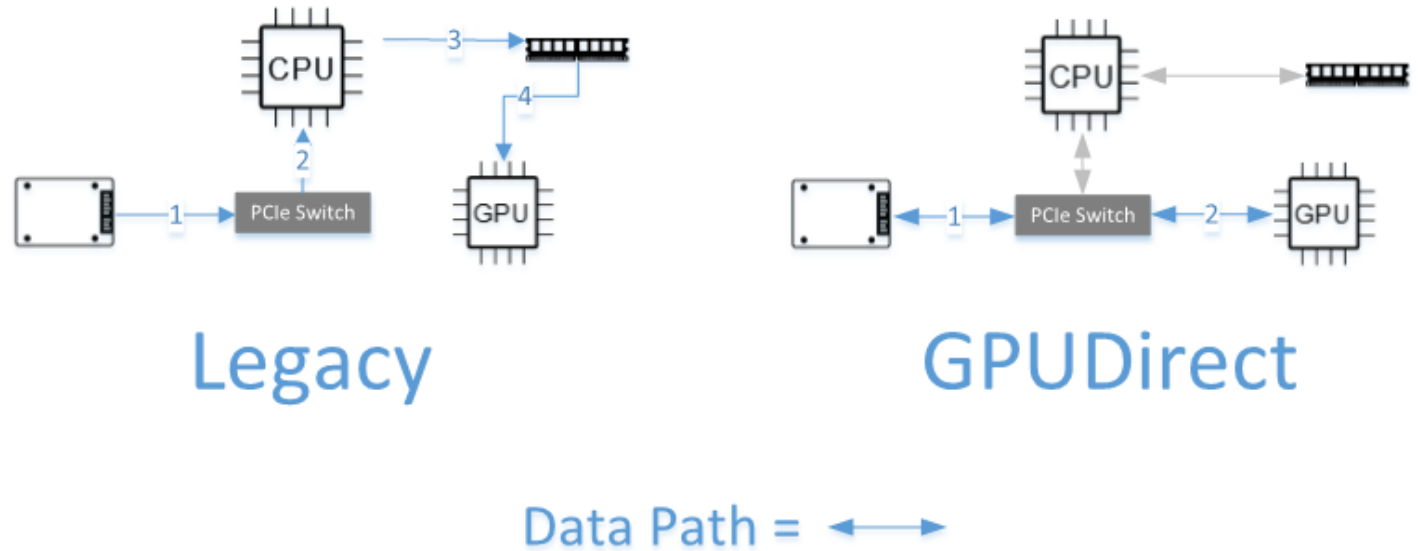
**The point here is that traditionally compute bound workloads are getting surprisingly close to being storage limited due to the architecture of AI training systems and the traditional data paths.**



# NVIDIA GPUDirect Storage

## What is GPUDirect Storage (GDS)

- GDS moves data directly between GPUs and Storage devices without using the CPU memory as a 'bounce buffer'
- Improves throughput and reduces latency



# GDS with Local NVMe Storage and V100s

---

## Test Configuration:

- SuperMicro SYS-4029GP-TVRT
- 2x Intel Xeon Platinum 8080M (28 Cores each)
- 3TB DRAM
- 8x Nvidia V100 SXM2 GPUs
- 8x NVMe SSDs

**Note: The NVMe SSDs are connected directly to CPUs, not on PCIe switches as in DGX2 or DGX A100**

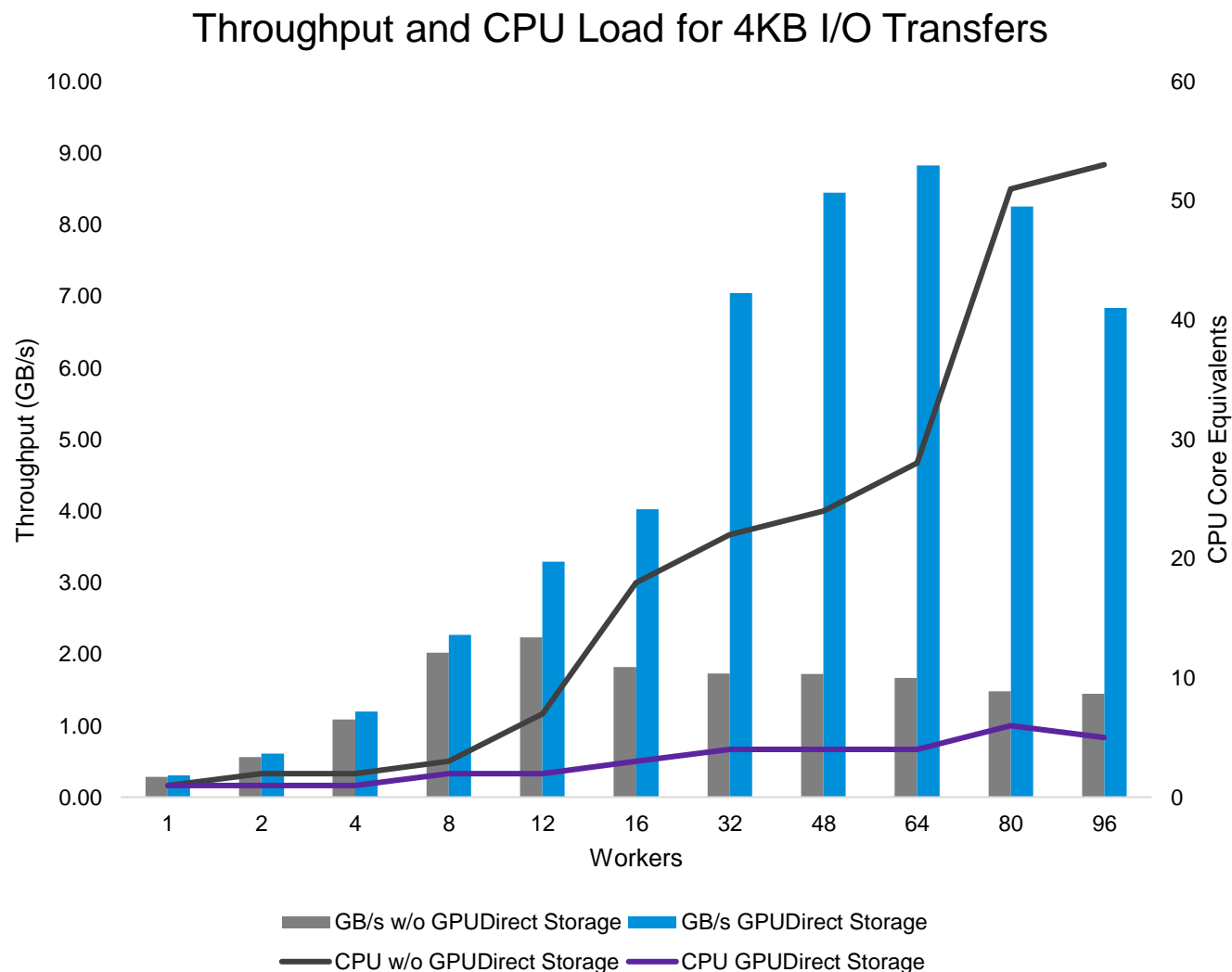




# CPU Utilization by Workers for 4k Transfers

## The Impact of Worker Count:

- Significant impact on Performance supplied by **GPUDirect Storage** vs. **CPU “bounce buffer”**.
- Higher worker count dramatically showing the GDS advantage increase.
- GDS Peak throughput at 64 workers while legacy path peaks at **12**.



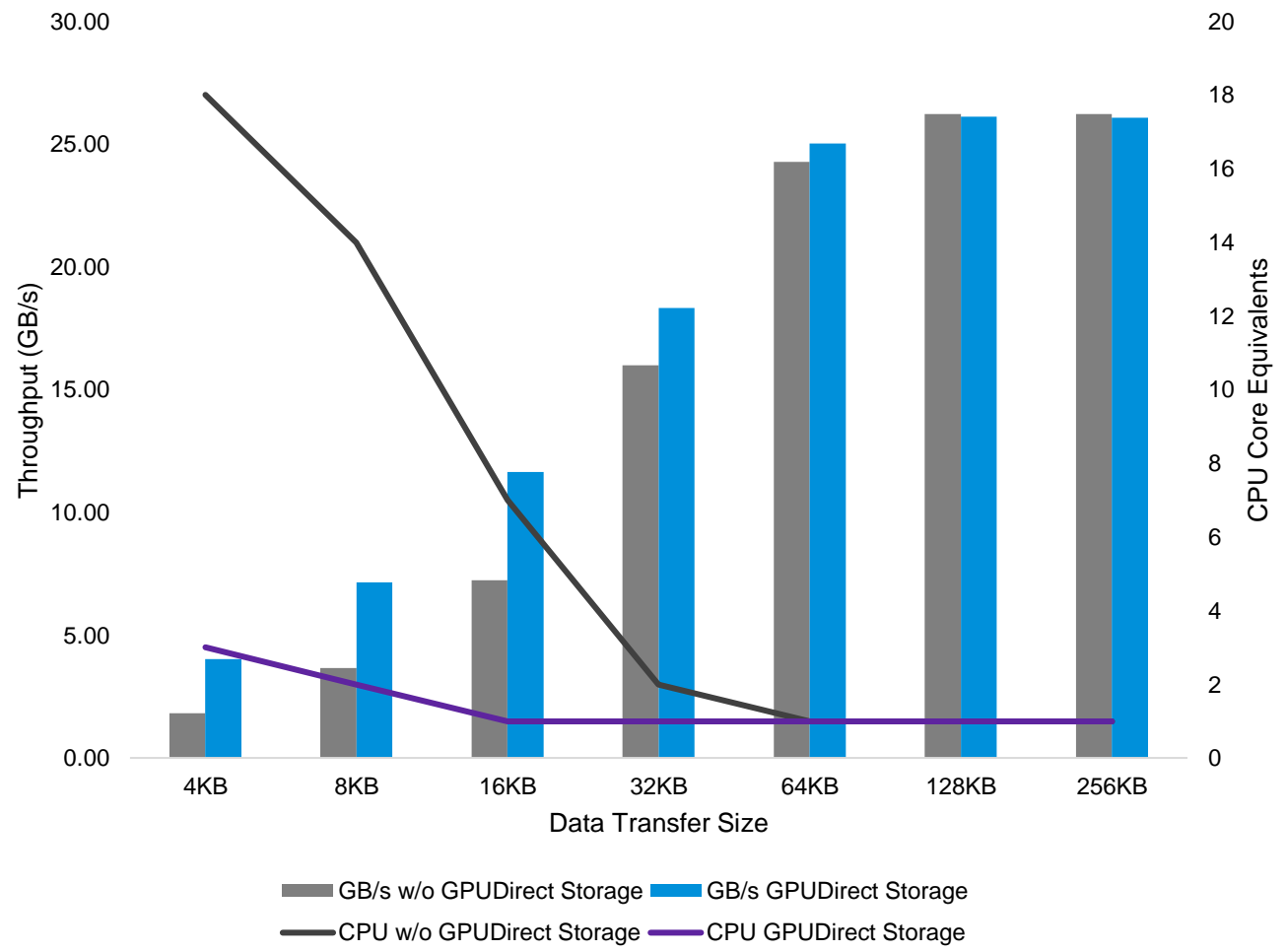
# CPU Utilization by IO Size

## Massive CPU savings for small IOs

Scaling the IO transfer size will mitigate the limitation of legacy data-path.

Throughput Improvement for Small to Medium IO transfer size.

CPU utilization and throughput by I/O transfer size for 8x GPUs by data path with 16 workers per GPU-NVMe pair



# GDS with Local NVMe Storage and V100s

- **16 Workers** per GPU is a “medium” load
- GDS consistently reduces latency for small to medium transfer size.

56% reduction at 4k

Negligible at large IOs

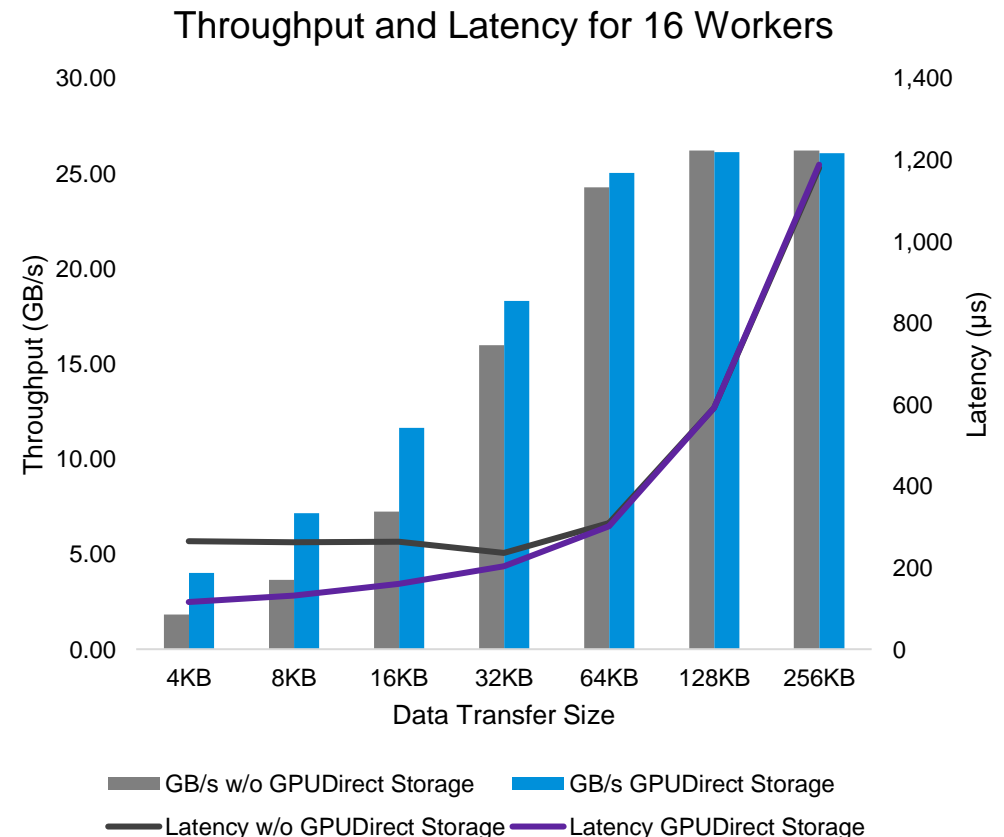
- GDS increases throughput at same load

121% increase at 4k

Negligible at large IOs

**The data here is interesting and shows that GDS can accelerate your storage in GPU systems**

Additional data from this testing presented in a blog post here:  
[Maximize Your Investment in Micron SSDs for AI/ML Workloads With NVIDIA GPUDirect Storage](#)



# Summary

- Storage clearly affect the training speed for AI applications.
- GPUDirect Storage (GDS)
  - Supplies considerable increases in total Throughput.
  - Latency decrease significantly.
  - Require lower CPU cores.

# Additional Collateral

- Micron Blog - [Maximize Your Investment in Micron SSDs for AI/ML Workloads With NVIDIA GPUDirect Storage](#)
- Nvidia Magnum-IO - <https://www.nvidia.com/en-us/data-center/magnum-io/>
- GTC 2020 Presentation by CJ Newburn – <https://www.nvidia.com/en-us/gtc/session-catalog/?ncid=so-face-85029#/session/1596756804120001kqY3>
- Webinar on demand (joint Micron-NVIDIA webinar): architecting to overcome AI challenges. <https://go.micron.com/Micron-NVIDIA-GDS-On-Demand-Registration.html>
- Joint Micron and NVIDIA podcast – [Overcoming AI data bottlenecks with NVIDIA GPUDirect Storage](#)
- NVIDIA Developer GDS page – <https://developer.nvidia.com/gpudirect-storage>
- Feel free to reach out:  
Or Lapid, Field Apps Engineer  
[olapid@micron.com](mailto:olapid@micron.com)  
+972-54-7716676



