

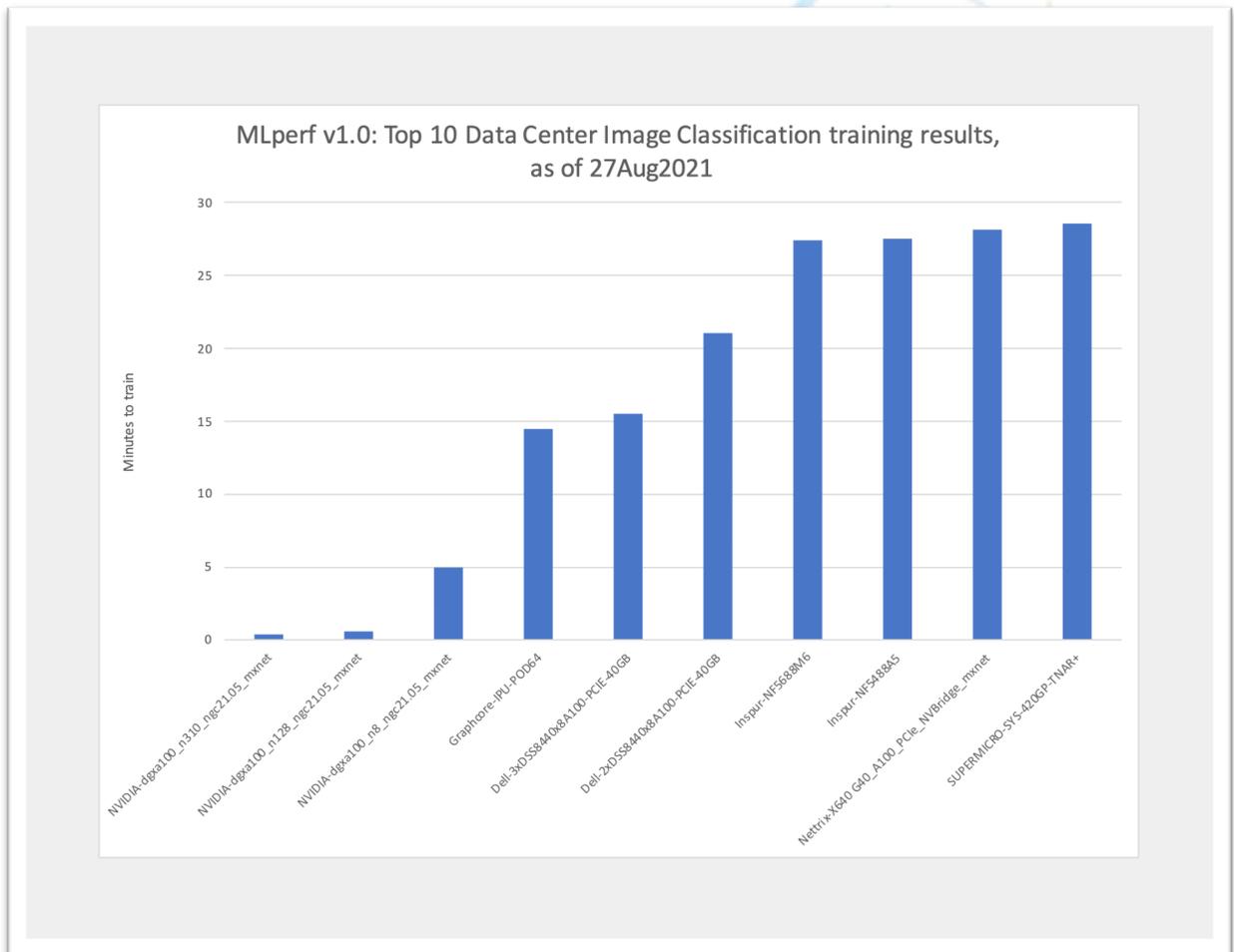
MLperf™ Performance Report

Silverton Consulting, Inc. StorInt™ Dispatch

This Storage Intelligence (StorInt™) dispatch covers the MLperf™ v1.0 series of AI-ML-DL model training and inferencing benchmarks¹. This report focuses on training activity for Data Center environments and workloads, which consists of image classification, image segmentation, object detection-light and -heavy, natural language processing, recommendation engine, and reinforcement learning model training, and runs them multiple times on a hardware configuration to a given accuracy level. The main MLperf training metric is the measured time to train a model, in minutes.

MLperf v1.0 Data Center training benchmark results

We start our discussion with image classification training benchmark in Figure 1.



¹ All MLperf inferencing and training results are available at <https://mlcommons.org/en/> as of 08/27/2021

Figure 1 Top 10 MLperf v1.0 DC Image Classification training results

In Figure 1, NVIDIA took the top 3 spots using 620 AMD EPYC 7742 CPUs/2480 NVIDIA A100-SXM4-80GB (400W) GPUs (#1), 256 EPYC 7742 CPUs/1024 A100-SXM4-80GB (400W) GPUs (#2) and 16 EPYC 7742 CPUs/64 A100-SXM4-80GB GPUs (#3), which took **0.4, 0.6 and 4.9 minutes to train**, respectively. One surprise is that the 16 CPU/64 GPU NVIDIA submission was only 8X slower than the 256 CPU/1024 GPU submission, by all rights it should have been 16X slower, with 16 less CPUs and GPUs. All these NVIDIA submissions ran Apache MXNet NVIDIA Release 21.05 software.

The Graphcore IPU-POD64 (#4) used 8 AMD EPYC 7742/64 GC200 IPU and trained in **14.5 minutes**, ~3X slower than the #3 NVIDIA system with 16 CPUs/64 A100 GPUs. Unclear whether one can compare an A100 GPU to a GC200 IPU but that seems logical. Given that, the IPU is ~3X slower than the A100 GPU. If we calculate the GC200 IPU memory properly, each IPU has 1.3TB (1472 IPU cores X 900MB of In-Processor-Memory™). This is ~17X the memory in an 80GB A100 GPU. So, a Graphcore system with ½ the CPUs and 17X more memory is 3X slower than a similar NVIDIA DGXA100 system in image classification training.

One consideration in their favor might be that of all the 10 submissions in Figure 1, Graphcore solution was the only one not using MXNet software, it used TensorFlow 1.15 (and Poplar SDK 2.1.0). Graphcore's Poplar SDK currently only supports PyTorch and Tensorflow 1 & 2.

All the remaining submissions shown in Figure 1 used A100 GPUs with different server hardware, some with 6 CPUs/24 A100-40GB GPUs (#5), 4 CPUs/16 A100-40GB GPUs (#6), and all the rest (#7-10) with 2 CPU/8 A100 GPUs and some GPUs with 40GB (#5,6 & 9) and others with 80GB (#7,8, & 10) of memory/GPU.

In Figure 2, we see similar results for NLP training.

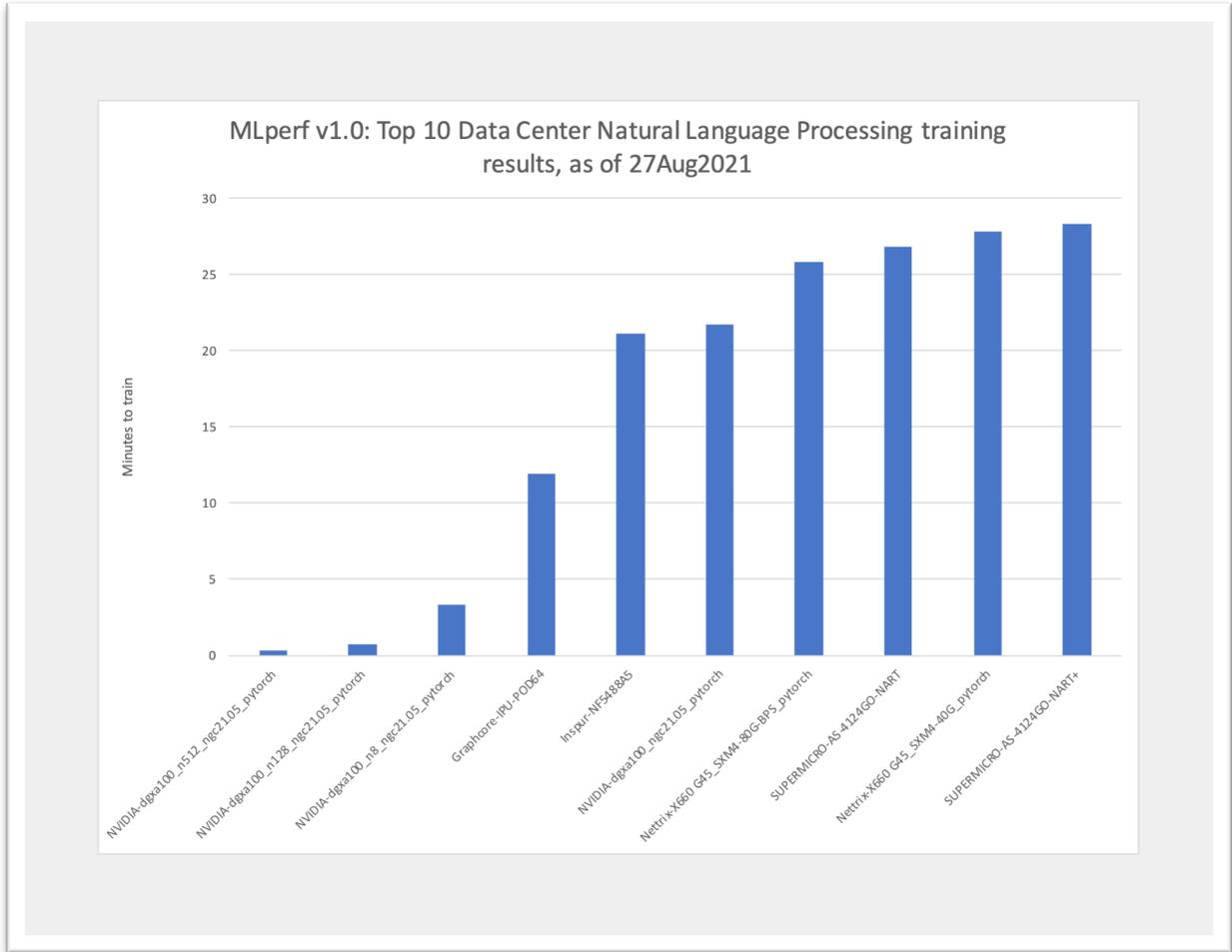


Figure 2 Top 10 MLperf v1.0 DC NLP training results

In Figure 2, the rankings are the exact same for #1-4 but differ after that. Again, the most interesting comparison here is the (#3) NVIDIA 16 CPUs/64 A100-80GB GPU which trained in **3.4 minutes** vs. the (#4) Graphcore 2 (believe this should be 8) CPUs/64 IPU which completed training in **12.0 minutes**. This result shows the Graphcore solution is ~3.5X slower than the NVIDIA A100 solution in NLP training.

For NLP training, the NVIDIA submissions used PyTorch NVIDIA Release 21.05 while Graphcore used PopART (and Poplar SDK 2.1.0), which we believe is Graphcore’s own DL framework.

One other item of note, but not shown in Figure 2, that is the Graphcore IPU-POD16 with 2 CPUs and 16 IPU came in at #11 with an NLP training time of **34.5 minutes**, which is ~3X slower than the IPU-POD64. Given that there is 4X the hardware. There seems to be something else slowing down the IPU-POD64. This submission also used the PopART software.

Similarly, Graphcore did submit an IPU-POD16 for the image classification benchmark that came in at #19 with an image classification training time of 37.1 minutes, or ~2.6X slower than the IPU-POD64 with 4X the hardware. This seems to re-iterate the story on the NLP benchmark. There's some slowdown with the 4 IPU-POD16s that go together to make up the IPU-POD64 system.

Significance

It's interesting to see that both Graphcore and NVIDIA submissions are not showing linear speedups with more hardware. The issue could be software or hardware bottlenecks, but with the amount of data being moved around in these benchmarks, we would lean towards hardware being the issue.

It's great to see NVIDIA dominate the benchmark results with Graphcore being an occasional exception, but where are the other GPUs. Both AMD and Intel have GPUs but we never see them in these benchmarks. Why?

This is our 5th MLperf performance report and we are starting to feel we understand them well. If there is something we missed or have an error in any of our analysis, please let us know and we would gladly fix it.

This report was sent out to subscribers as part our **free, monthly Storage Intelligence e-newsletter**. If you are interested in receiving future storage performance analyses along with recent product announcement summaries, please use the QR code (below right) to sign up for your own copy.

Silverton Consulting, Inc., is a U.S.-based Storage, Strategy & Systems consulting firm offering products and services to the data storage community

Newsletter signup QRcode

