Home (/) > Storage (/t/storage/) > PowerFlex (/t/powerflex-14/) > Blogs (/t/blogs-44/) > Can I do that AI thing on Dell PowerFlex?

# Can I do that AI thing on Dell PowerFlex?

Thu, 20 Jul 2023 21:08:09 -0000 | Read Time: 0 minutes

Tony Foster    Tony Foster

The simple answer is Yes, you can do that AI thing with Dell PowerFlex. For those who might have been busy with other things, AI stands for Artificial Intelligence and is based on trained models that allow a computer to "think" in ways machines haven't been able to do in the past. These trained models (neural networks) are essentially a long set of IF statements (layers) stacked on one another, and each IF has a 'weight'. Once something has worked through a neural network, the weights provide a probability about the object. So, the AI system can be 95% sure that it's looking at a bowl of soup or a major sporting event. That, at least, is my

Most recently, AI has been made famous by large language models (LLMs) for conversational AI applications like ChatGPT. Though these applications have stoked fears that AI will take over the world and destroy humanity, that has yet to be seen. Computers still can do only what we humans tell them to do, even LLMs, and that means if something goes wrong, we their creators are ultimately to blame. (See 'Godfather of AI' leaves Google, warns of tech's dangers (https://apnews.com/article/ai-godfather-google-geoffery-hinton-fa98c6a6fddab1d7c27560f6fcbad0ad).)

The reality is that most organizations aren't building world destroying LLMs, they are building systems to ensure that every pizza made in their factory has exactly 12 slices of pepperoni evenly distributed on top of the pizza. Or maybe they are looking at loss prevention, or better traffic light timing, or they just want a better technical support phone menu. All of these are uses for AI and each one is constructed differently (they use different types of neural networks).

We won't delve into these use cases in this blog because we need to start with the underlying infrastructure that makes all those ideas "AI possibilities." We are going to start with the infrastructure and what many now consider a basic (by today's standards) image classifier known as ResNet-50 v1.5. (See ResNet-50: The Basics and a Quick Tutorial (https://datagen.tech/guides/computer-vision/resnet-50/).)
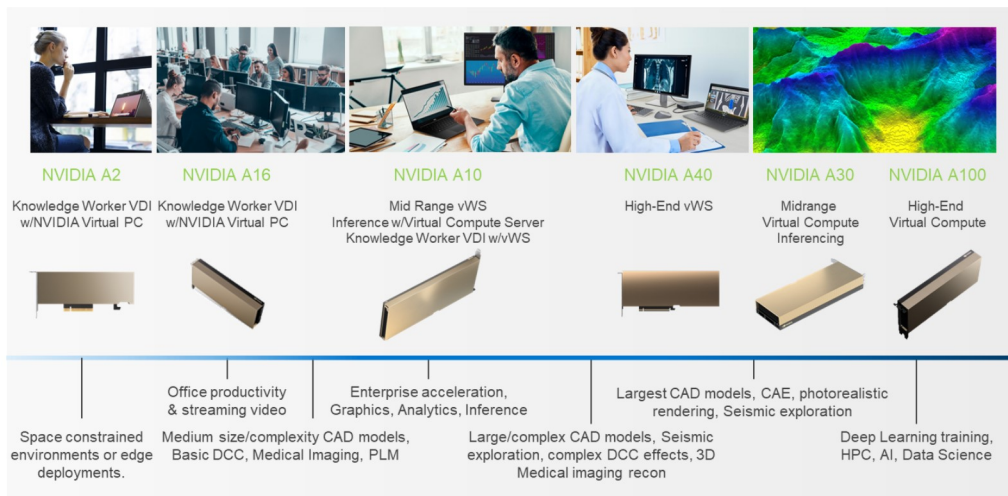
That's also what the PowerFlex Solution Engineering team did in the validated design (https://infohub.delltechnologies.com/t/dell-validated-design-for-virtual-gpu-with-vmware-and-nvidia-on-powerflex/) they recently published. This design details the use of ResNet-50 v1.5  in a VMware vSphere environment using NVIDIA AI Enterprise (https://www.nvidia.com/en-us/data-center/products/ai-enterprise/) as part of PowerFlex environment. They started out with the basics of how a virtualized NVIDIA GPU works well in a PowerFlex environment. That's what we'll explore in this blog – getting started with AI workloads, and not how you build the next AI supercomputer (though you could do that with PowerFlex as well).

In this validated design, they use the NVIDIA A100 (PCIe) GPU (https://www.nvidia.com/en-us/data-center/a100/) and virtualized it in VMware

us/technologies/multi-instance-gpu/) technology.

NVIDIA's MIG technology allows administrators to partition a GPU into a maximum of seven GPU instances. Being able to do this provides greater control of GPU resources, ensuring that large and small workloads get the appropriate amount of GPU resources they need without wasting any.

PowerFlex supports a large range of NVIDIA GPUs for workloads, from VDI (Virtual Desktops) to high end virtual compute workloads like AI. You can see this in the following diagram where there are solutions for "space constrained" and "edge" environments, all the way to GPUs used for large inferencing models. In the table below the diagram, you can see which GPUs are supported in each type of PowerFlex node. This provides a tremendous amount of flexibility depending on your workloads.
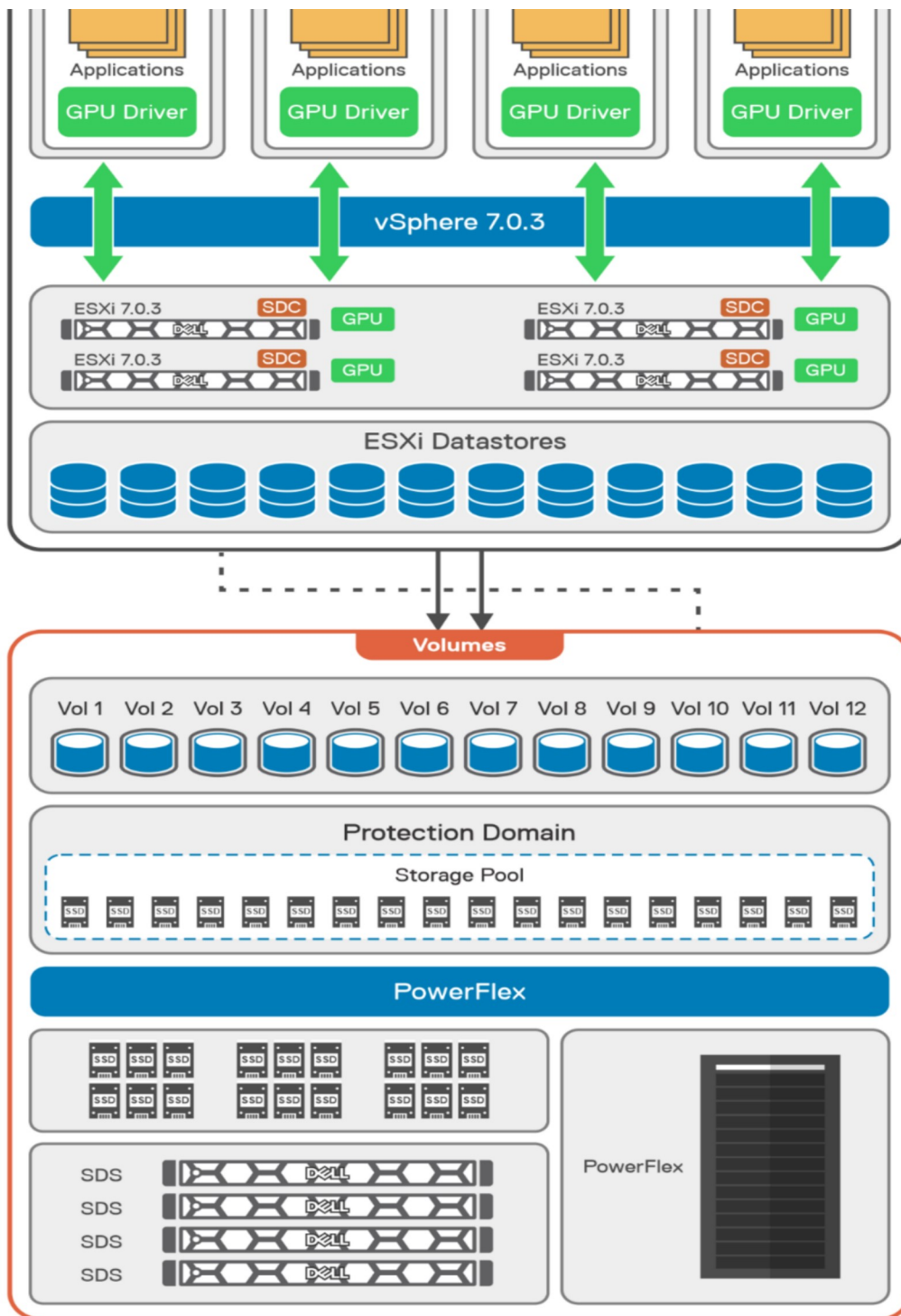
| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | A40 | (x) | (x) | (x) | ✓ | (x) | ✓ | ✓ |
| | A30 | (x) | (x) | (x) | ✓ | (x) | ✓ | ✓ |
| | A10 | (x) | (x) | (x) | ✓ | (x) | ✓ | ✓ |
| | A16 | (x) | (x) | (x) | ✓ | (x) | ✓ | ✓ |
| | A2 | ✓ | ✓ | ✓ | ✓ | (x) | ✓ | ✓ |
| | T4 | ✓ | ✓ | ✓ | ✓ | (x) | ✓ | ✓ |

The validated design describes the steps to configure the architecture and provides detailed links to the NVIDIAand VMware documentation for configuring the vGPUs, and the licensing process for NVIDIA AI Enterprise.

These are key steps when building an AI environment. I know from my experience working with various organizations, and from teaching, that many are not used to working with vGPUs in Linux. This is slowly changing in the industry. If you haven't spent a lot of time working with vGPUs in Linux, be sure to pay attention to the details provided in the guide (https://infohub.delltechnologies.com/t/dell-validated-design-for-virtual-gpu-with-vmware-and-nvidia-on-powerflex/). It is important and can make a big difference in your performance.
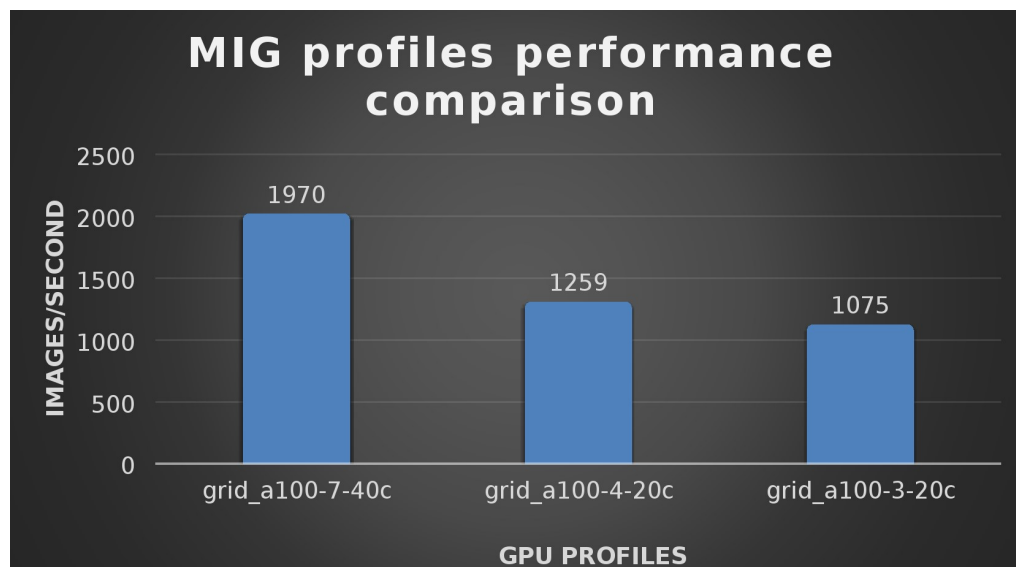
The following diagram shows the validated design's logical architecture. At the top of the diagram, you can see four Ubuntu 22.04 Linux VMs with the NVIDIA vGPU driver loaded in them. They are running on PowerFlex hosts with VMware ESXi deployed. Each VM contains one NVIDIA A100 GPU configured for MIG operations. This configuration leverages a two-tier architecture where storage is provided by separate PowerFlex software defined storage (SDS) nodes.
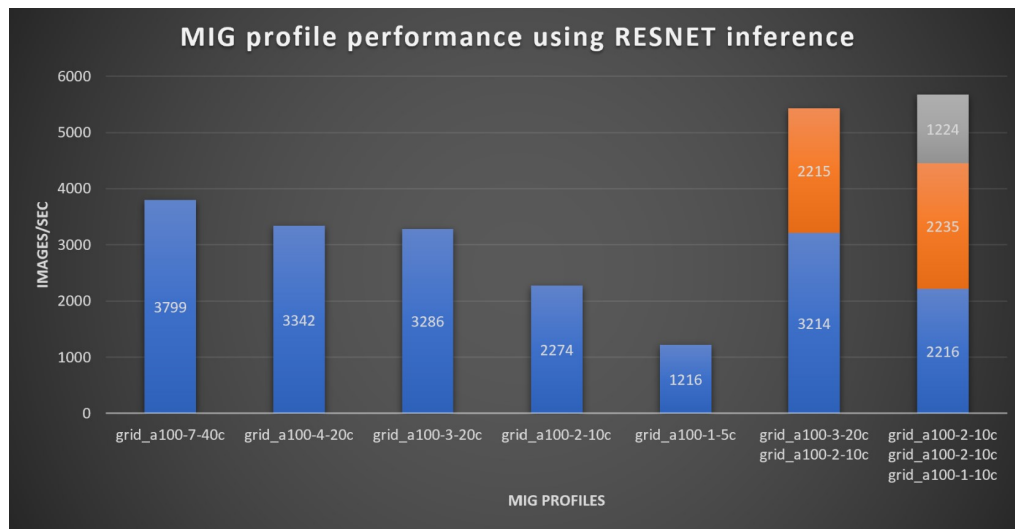
A design like this allows for independent scalability for your workloads. What I mean by this is during the training phase of a model, significant storage may be required for the training data, but once the model clears validation and goes into production, storage requirements may be drastically different. With PowerFlex you

team validated it using ResNet-50 v1.5 using the ImageNet (http://www.image-net.org/) 1K data set. For this validation they enabled several ResNet-50 v1.5 TensorFlow features. These include Multi-GPU training with Horovod (https://github.com/horovod/horovod), NVIDIA DALI (https://docs.nvidia.com /deeplearning/dali/release-notes/index.html), and Automatic Mixed Precision (AMP) (https://developer.nvidia.com/automatic-mixed-precision). These help to enable various capabilities in the ResNet-50 v1.5 model that are present in the environment. The paper then describes how to set up and configure ResNet-50 v1.5, the features mentioned above, and details about downloading the ImageNet data.

At this stage they were able to train the ResNet-50 v1.5 deployment. The first iteration of training used the NVIDIA A100-7-40C vGPU profile (https://docs.nvidia.com/grid/13.0/grid-vgpu-user-guide/index.html#vgpu-types-nvidia-a100-pcie-40gb). They then repeated testing with the A100-4-20C vGPU profile and the A100-3-20C vGPU profile. You might be wondering about the A100-2-10C vGPU profile and the A100-1-5C profile. Although those vGPU profiles are available, they are more suited for inferencing, so they were not tested.



The results from validating the training workloads for each vGPU profile is shown in the following graph. The vGPUs were running near 98% capacity according to nvitop (https://pypi.org/project/nvitop/) during each test. The CPU utilization was 14% and there was no bottle neck with the storage during the tests.

It's worth noting that the last two columns show the inferencing running across multiple VMs, on the same ESXi host, that are leveraging MIG profiles. This also shows that GPU resources are partitioned with MIG and that resources can be precisely controlled, allowing multiple types of jobs to run on the same GPU without impacting other running jobs.

This opens the opportunity for organizations to align consumption of vGPU resources in virtual environments. Said a different way, it allows IT to provide "show back" of infrastructure usage in the organization. So if a department only needs an inferencing vGPU profile, that's what they get, no more, no less.

It's also worth noting that the results from the vGPU utilization were at 88% and CPU utilization was 11% during the inference testing.

These validations show that a Dell PowerFlex environment can support the foundational components of modern-day AI. It also shows the value of NVIDIA's MIG technology to organizations of all sizes: allowing them to gain operational efficiencies in the data center and enable access to AI.

Which again answers the question of this blog, can I do that AI thing on Dell PowerFlex… Yes you can run that AI thing! If you would like to find out more about how to run your AI thing on PowerFlex, be sure to reach out to your Dell representative.

# Resources

fa98c6a6fddab1d7c27560f6fcbad0ad)

- ResNet-50: The Basics and a Quick Tutorial (https://datagen.tech/guides /computer-vision/resnet-50/)
- Dell Validated Design for Virtual GPU with VMware and NVIDIA on PowerFlex (https://infohub.delltechnologies.com/t/dell-validated-design-for-virtual-gpu-with-vmware-and-nvidia-on-powerflex)
- NVIDIA NGC Catalog ResNet v1.5 for PyTorch (https://catalog.ngc.nvidia.com /orgs/nvidia/resources/resnet_50_v1_5_for_pytorch)
- NVIDIA AI Enterprise (https://www.nvidia.com/en-us/data-center/products/ai-enterprise/)
- NVIDIA A100 (PCIe) GPU (https://www.nvidia.com/en-us/data-center/a100/)
- NVIDIA Virtual GPU Software Documentation (https://docs.nvidia.com /grid/latest/index.html)
- NVIDIA A100-7-40C vGPU profile (https://docs.nvidia.com/grid/13.0/grid-vgpu-user-guide/index.html#vgpu-types-nvidia-a100-pcie-40gb)
- NVIDIA Multi-Instance GPU (MIG) (https://www.nvidia.com/en-us/technologies /multi-instance-gpu/)
- NVIDIA Multi-Instance GPU User Guide (https://docs.nvidia.com/datacenter /tesla/mig-user-guide/index.html)
- Horovod (https://github.com/horovod/horovod)
- ImageNet (http://www.image-net.org/)
- DALI (https://docs.nvidia.com/deeplearning/dali/release-notes/index.html)
- Automatic Mixed Precision (AMP) (https://developer.nvidia.com/automatic-mixed-precision)
- nvitop (https://pypi.org/project/nvitop/)

**Author**: Tony Foster

Sr. Principal Technical Marketing Engineer

| Twitter: | @wonder_nerd (https://twitter.com/wonder_nerd) |
| --- | --- |
| LinkedIn: | https://linkedin.com/in/wondernerd/ (https://linkedin.com/in/wondernerd/) |
| Personal Blog: | https://wondernerd.net (https://wondernerd.net) |
| Location: | The Land of Oz [-6 GMT] |

# Related Blog Posts



(/p/accelerating-distributed-training-in-a-multinode-virtualized-environment/)

AI    deep learning    NVIDIA    PowerEdge    VMware    GPU    PowerScale

## Accelerating Distributed Training in a Multinode Virtualized Environment

(/p/accelerating-distributed-training-in-a-multinode-virtualized-environment/)
Srinivas Varadharajan Bala Chandrasekaran Prem Pradeep Motgi Sarvani Vemulapalli …

Fri, 13 May 2022 13:57:13 -0000 | Read Time: 0 minutes
(/p/accelerating-distributed-training-in-a-multinode-virtualized-environment/)
Introduction
In the age of deep learning (DL), with complex models, it is vital to have a system that allows
faster distributed training. Depending on the application, some DL models require more



(/p/comparison-of-top-accelerators-from-dell-technologies-mlperf-tm-inference-v3-0-submission/)

deep learning    NVIDIA    PowerEdge    machine learning    GPU    MLPerf

# Comparison of Top Accelerators from Dell Technologies' MLPerf™ Inference v3.0 Submission

(/p/comparison-of-top-accelerators-from-dell-technologies-mlperf-tm-inference-v3-0-submission/)

Manpreet Sokhi Frank Han Rakshith Vasudev   Manpreet Sokhi,  Frank Han,  Rakshith Vasudev

Fri, 21 Apr 2023 21:43:39 -0000 | Read Time: 0 minutes
(/p/comparison-of-top-accelerators-from-dell-technologies-mlperf-tm-inference-v3-0-submission/)
Abstract

Dell Technologies recently submitted results to MLPerf™ Inference v3.0 in the closed division.

logo