

Technical Review

Accelerating the Artificial Intelligence Journey with Dell EMC Ready Solutions for AI

Date: August 2018 Author: Jack Poller, Senior Analyst

Abstract

This ESG Technical Review documents evaluation of Dell EMC Ready Solutions for AI. We focused on understanding the performance and ease of use of the Ready Solutions for AI with optimized designs for machine learning and deep learning. To validate the full stack performance, we measured the number of images per second processed when training the AlexNet and ResNet50 deep learning networks and evaluated how the integrated solutions can simplify and accelerate AI deployment. The Deep Learning with NVIDIA design featuring Isilon significantly outperformed the competition in training time, delivering 2.9 times the performance of one competitor for an AlexNet deep learning neural network in a GPU accelerated environment and 2.3 times another competitor for a ResNet50 deep learning neural network.

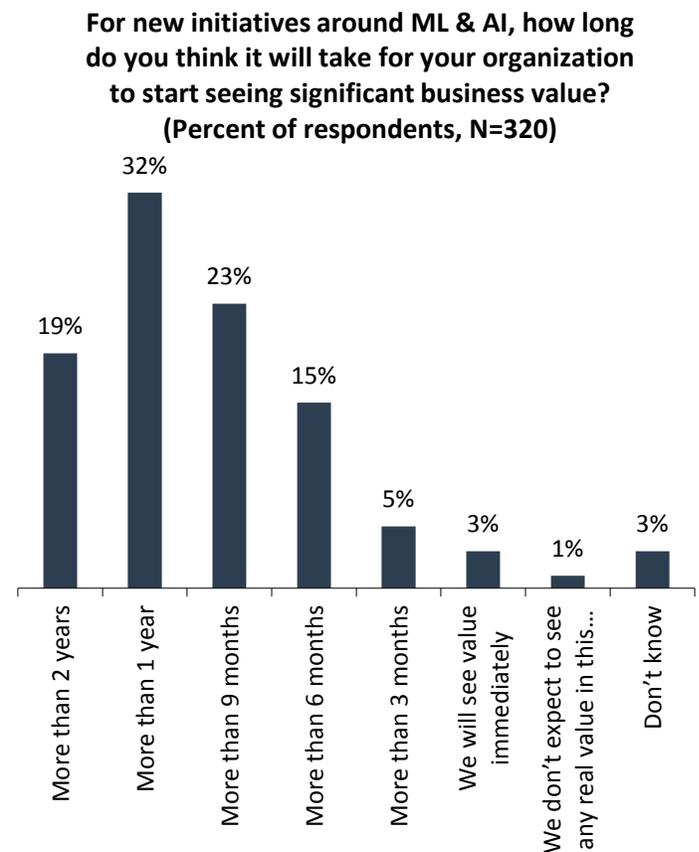
The Challenges

A single graphics processor (GPU) can now achieve 100 teraFLOPS¹ due to advances in microprocessor architecture and design, such as the ability to combine hundreds or thousands of processor cores into a single chip. This massive amount of computing power accelerates training of massively parallel and iterative artificial intelligence (AI) models, helping deep and machine learning to become viable techniques for any business to extract additional value from its data.

Machine learning and deep learning routinely leverage unstructured data such as images, video, and streaming sensor data—data that is often uncompressible and quickly scales from tens of TBs to tens of PBs. Organizations are challenged to develop AI solutions that can both manage data at scale and provide the network and storage performance to feed the data-hungry, massively concurrent compute layer.

However, AI is in its infancy, and lacks a standardized infrastructure stack. Thus, it takes a significant amount of time for organizations to develop their AI infrastructures and obtain results that impact their business. According to recent ESG research, 19% of organizations don't expect to extract business value from their AI efforts for two years. Another 32% of organizations think it will take more than a year, and 23% think it will take more than nine months (see Figure 1).²

Figure 1. Time to Value Expectations for AI



Source: Enterprise Strategy Group

¹ 1 teraFLOPS is one million million (10¹²) floating-point operations per second.

² Source: ESG Survey, *Machine Learning and Artificial Intelligence Trends*, June 2017. All ESG research references and charts in this technical review have been taken from this survey unless otherwise noted.

Dell EMC Ready Solutions for AI

With Ready Solutions for AI, Dell EMC has created standardized infrastructure stacks for machine learning (ML) and deep learning (DL) to accelerate the time to business value.

Two Dell EMC Ready Solutions for AI are available today:

- **Machine Learning with Hadoop**—Optimized for machine learning and deep learning with Hadoop, it includes:
 - Cloudera Data Science Cluster: head node and two worker nodes with 960GB-1.92TB direct attached SSD storage.
 - Hadoop nodes: starting with three infrastructure nodes and seven worker nodes, scaling out to thousands of nodes.
 - 25GbE Ethernet networking with Dell EMC Open Network Switches.
 - Software stack: Cloudera Manager, Cloudera Data Science Workbench, Cloudera Enterprise Data Hub, Spark, and Dell EMC Data Science Provisioning Engine.
 - Frameworks/libraries: BigDL.
- **Deep Learning with NVIDIA**—Optimized for deep learning with GPU acceleration, it includes:
 - 1 PowerEdge 740 Head node: dual-processor head node with 12 x 10TB direct attached SAS drives.
 - 4 PowerEdge C4140 Worker nodes: dual-processor nodes with 384GB memory and up to four NVIDIA Tesla V100 GPUs each with the ability to scale out to thousands of nodes.
 - 100Gb/s networking with Mellanox Infiniband switches, and Dell EMC Open Network top-of-rack switches.
 - Storage: Isilon F800 All-flash Scale-out NAS, options for 96, 192, or 924TB flash capacity per chassis, 15GB/s bandwidth per chassis, 8x 40GbE networking per chassis. Scales out up to 33 PB and up to 540GB/s bandwidth per cluster.
 - Software stack: Bright Cluster Manager for Data Science, and Dell EMC Data Science Provisioning Portal.
 - Frameworks/libraries: Caffe 2, MXNET, TensorFlow, NVIDIA CUDA Deep Neural Network library (cuDNN), and NVIDIA CUDA basic linear algebra subroutines (cuBLAS).



Dell EMC Ready Solutions for AI come with deployment services to accelerate time to results and single contact support for the complete hardware and software stack.

These validated hardware and software stacks combine Dell EMC PowerEdge servers, Dell EMC Isilon storage, NVIDIA GPUs, high-speed networking, data science software, and AI libraries and frameworks into preconfigured, scalable, tuned systems. Organizations deploying Ready Solutions for AI benefit from:

- **Fast deployment**—Rather than forcing the organization to select, configure, integrate, and tune components into an AI stack, Dell EMC Ready Solutions for AI are validated systems deployed by Dell EMC services, shrinking the time to deploy an AI environment from months to weeks while reducing skillset requirements and operational risk.
- **Simplified configuration**—Both designs increase data scientist productivity by offering self-service access to resources for machine learning and deep learning including frameworks and libraries such as BigDL, TensorFlow, Caffe, Neon, cuDNN, and cuBLAS. The Deep Learning with NVIDIA design includes Dell EMC Data Science Provisioning Portal, which reduces the steps it takes to configure a data scientist's workspace to just five clicks. Machine Learning with Hadoop includes Cloudera's Data Science Workbench and Dell EMC Data Science Engines—containers that work with Data Science Workbench to configure the BigDL framework.
- **Simplified IT operations**—Each design includes a single console for monitoring the health and configuration of the cluster. Deep Learning with NVIDIA includes Bright Computing's Bright Cluster Manager that offers integrations with

Dell Remote Access Controller for PowerEdge servers to monitor and manage the health and configuration of the cluster. Machine Learning with Hadoop includes Cloudera Manager for monitoring and configuration management of the Hadoop cluster.

- **Rapid scalability**—Dell EMC designed Ready Solutions for AI for rapid scalability. Organizations can increase compute power by adding compute nodes to the cluster with just a few mouse clicks. Storage can be scaled out by non-disruptively adding additional nodes, which linearly increases storage performance.

Understanding the Opportunities and Challenges of AI

While artificial intelligence dates to the early days of the computer age, practical and achievable machine learning and deep learning are relatively new fields, and there is a general lack of expertise and guidance available. Obtaining meaningful results that impact business outcomes requires massive computing power to process equally massive data sets using complex software frameworks and libraries.

Creating an AI infrastructure stack requires both AI expertise to assemble the proper combinations of software solutions and systems and integration expertise to assemble and tune the proper combinations of hardware solutions to create an efficient, scalable, and cost-effective system.

IT staff and data scientists must work in concert to select and acquire compute servers, GPUs, storage, and network. Once all systems arrive, are physically installed, and are powered on, IT must install, configure, and test the storage, networking, and operating systems. Next, the IT and/or data science teams need to install, configure, test, and tune the selected configurations of open source AI frameworks, libraries, and orchestration software. Finally, the data scientists need to validate the AI system. After this long and drawn-out process, which can take months, data scientists can start to create AI models. And minimal changes in the stack can lead to mediocre performance or even failure.

The major public cloud providers offer GPU-accelerated AI compute instances and AI libraries, enabling organizations to jump-start their AI programs. However, the public cloud offerings lack reference configurations, customer solution centers, and consulting, forcing data scientists to learn by themselves how to best configure and tune their AI stack. In addition, data locality and data movement between the cloud, the edge, and the core can impact both performance and costs, frequently making an on-premises solution a better choice.

AI models yield better results with larger data sets, and data scientists often analyze terabytes to petabytes of data. Organizations using the public cloud must pay for CPU time, GPU time, data storage, data ingress (network cost to transfer data into the public cloud), and recurring inferencing charges. While using the public cloud transforms capital expenditures into operating expenditures, costs are extremely variable and may not be predictable; when AI models don't converge, organizations can be surprised by monthly bills orders of magnitude greater than expected.

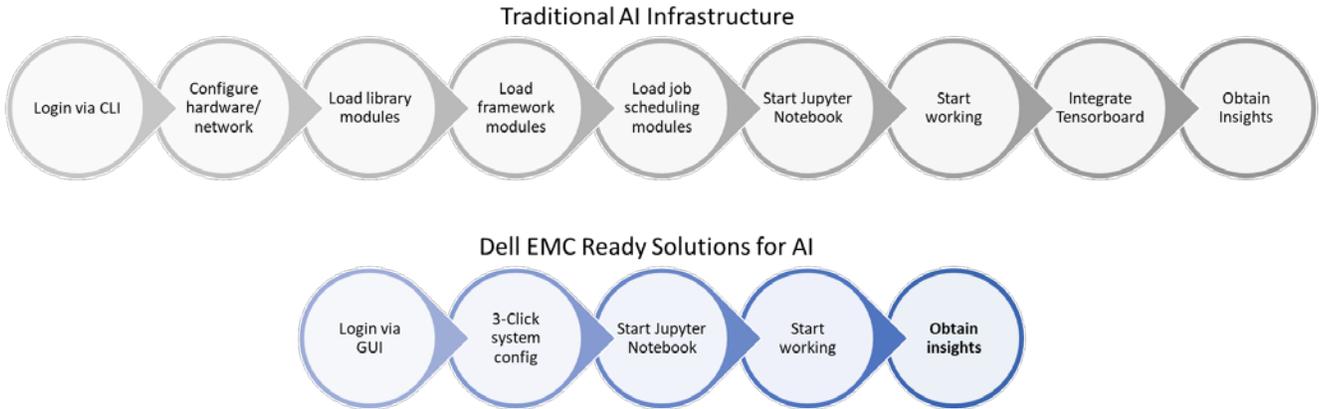
Simplifying AI Deployment

ESG started by evaluating how the Dell EMC Ready Solutions for AI simplify deployment of the AI infrastructure stack and accelerate time to results for the data scientist. Dell EMC Ready Solutions for AI come with all necessary software, compute, storage, and networking hardware that is installed onsite by Dell EMC professional services.

IT and data scientists can skip the time consuming and convoluted effort of installing and configuring operating systems, AI libraries, orchestration, and management software, saving weeks to months of effort.

The solutions include a self-service environment for data scientists to obtain cluster resources and configure frameworks and libraries for their work. These GUI systems simplify the effort for data scientists and IT to configure their workspace and manage the cluster. Whereas traditionally, data scientists would use the command line to configure their environment, these GUIs automate and orchestrate many tasks, enabling scientists to administer clusters as a single entity; provision hardware, operating system, and software; manage the cluster operation; provision workloads; and obtain results. For example, Deep Learning with NVIDIA includes the Data Science Provisioning Portal.

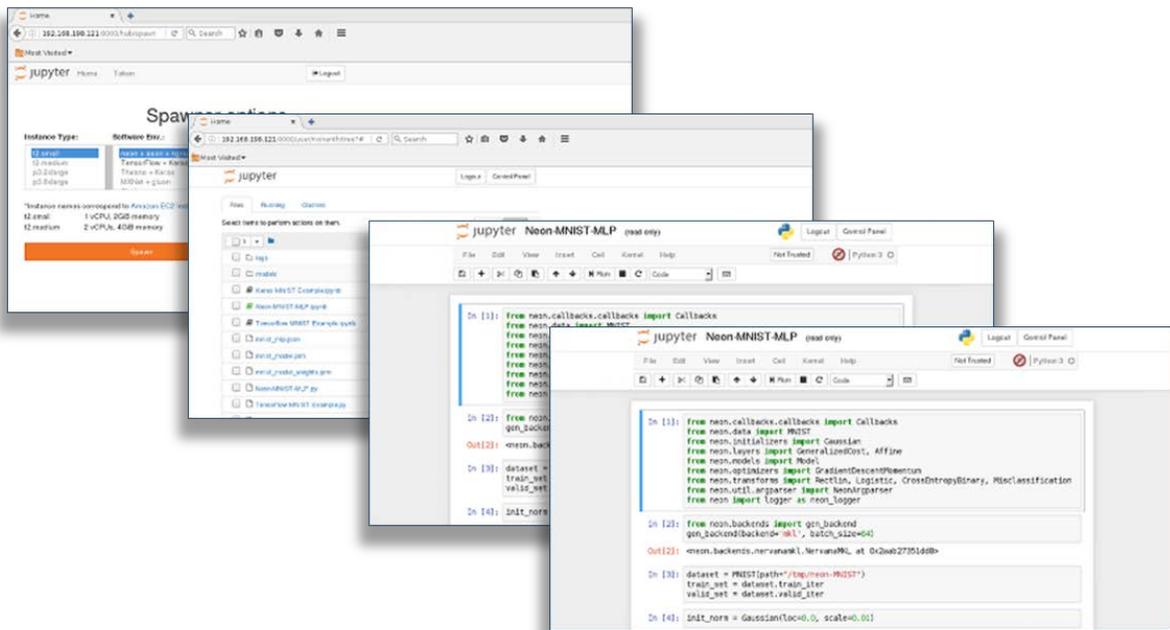
Figure 2. Level of Effort for Traditional AI Infrastructure and Dell EMC Ready Solutions for AI



Source: Enterprise Strategy Group

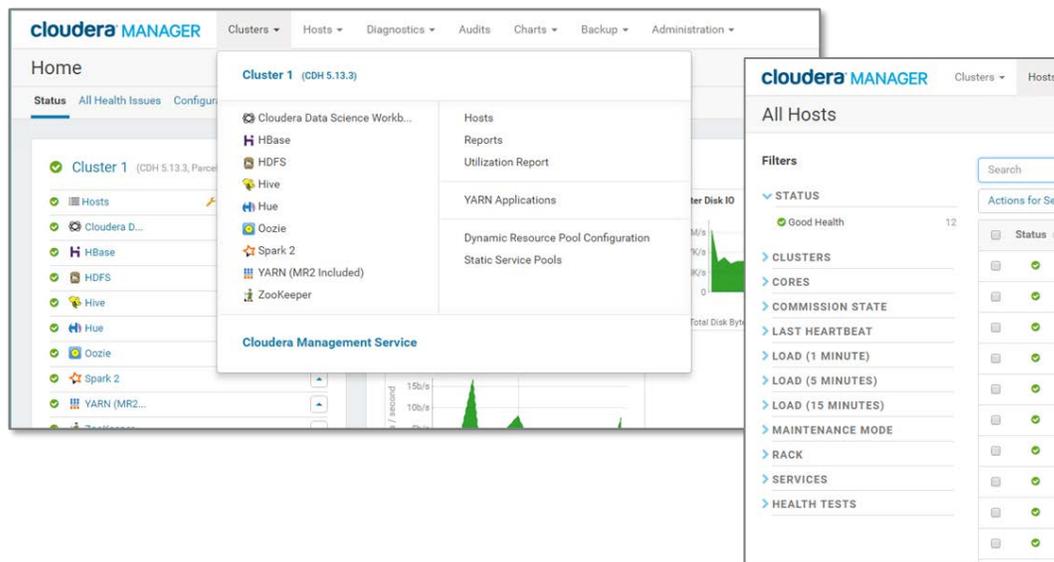
ESG logged in to the Data Science Provisioning Portal GUI. As shown in Figure 3, the Dell EMC Data Science Provisioning Portal required just three mouse clicks to select the compute and storage resources, the library modules, and the framework modules. Rather than using command lines, we were able to train our AI models and obtain insights and results from within the GUI.

Figure 3. Dell EMC Data Science Provisioning Portal



Source: Enterprise Strategy Group

We also reviewed the Cloudera System Manager included with Machine Learning with Hadoop, as shown in Figure 4. The dashboard view displayed the status, throughput, and load for each cluster and the cluster’s components. Using the pulldown menus, we could select and manage the entire cluster as a single entity or manage individual cluster components.

Figure 4. Cludera System Manager

Source: Enterprise Strategy Group

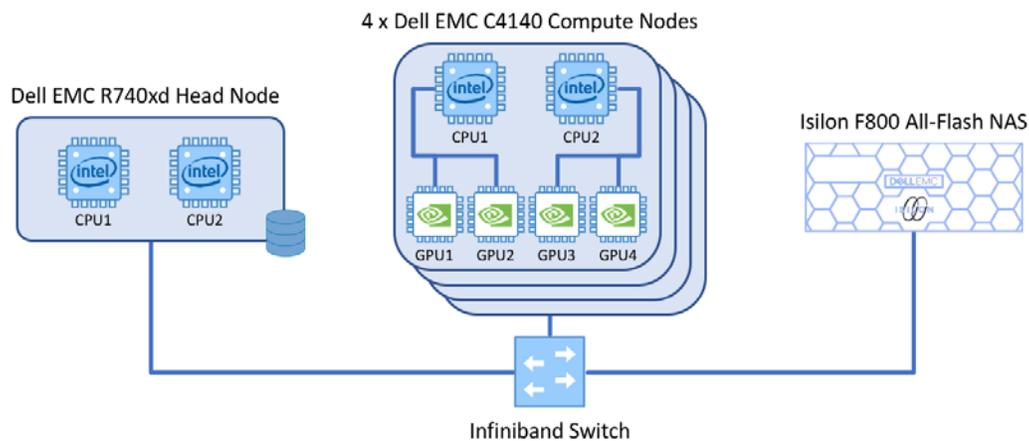
Why This Matters

Machine learning and deep learning are immature technologies, and the complete AI infrastructure stack is complex, requiring the integration of a variety of hardware and software components from many vendors, an intricate process that can take many months. Selecting the wrong components or misconfiguring the integration can cause I/O bottlenecks, which lead to poor performance, and system errors, which yield miserable results, limiting ROI from AI investments.

ESG validated that the Dell EMC Ready Solutions for AI provided a complete, integrated solution with CPUs, GPUs, networking, and scale-out storage. After installation by Dell EMC, data scientists can go from power-on to evaluating AI models in just a few mouse clicks using the included system management software. Rather than spend time working with IT to select and acquire components; configure the network; or install and configure operating systems, libraries, and frameworks; data scientists can proceed immediately to creating AI solutions, simplifying and shortening deployment time from months to weeks.

Accelerating AI Model Development

ESG evaluated how the Dell EMC Ready Solutions for AI accelerated the model training processes. We started with an environment consisting of the Deep Learning with NVIDIA system, as shown in Figure 5. The solution consisted of a five-server compute complex. One server, designated the head node, was used for system management, and the remaining four servers, each with two Intel Xeon Gold 6148 20-core processors, 384GB RAM, and four NVIDIA Tesla V100 GPUs, were used as compute nodes. As tested, the system included a single chassis of Isilon F800 All-Flash Scale-out NAS with 15GB/sec bandwidth and 192TB capacity. All servers were connected through Mellanox switches using 100Gb/s InfiniBand, and the Isilon was connected with eight 40GbE Ethernet links.

Figure 5. Ready Solutions for Deep Learning Test Bench

Source: Enterprise Strategy Group

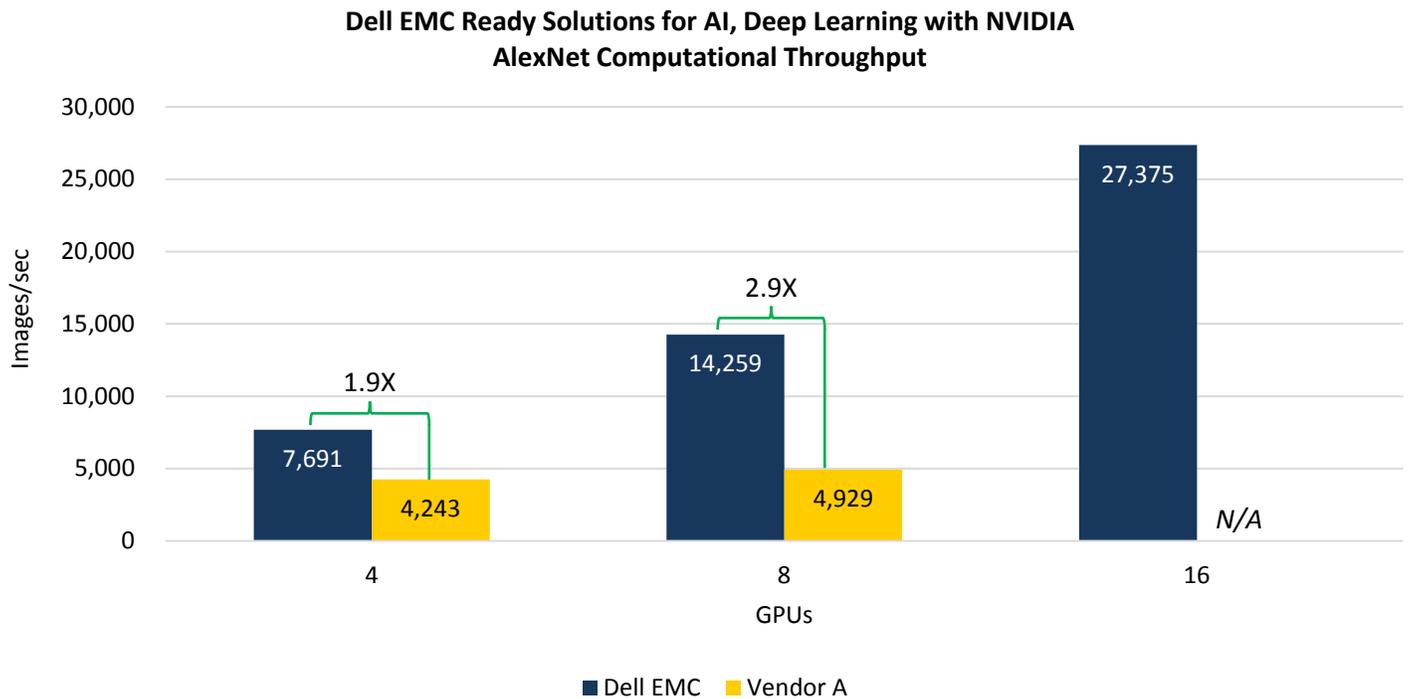
ESG used two different industry-standard benchmarks for GPU-accelerated infrastructure stacks to characterize the performance of Dell EMC Ready Solution for AI, Deep Learning with NVIDIA. We started with the AlexNet, an image classifier that can classify images into 1,000 object categories such as keyboard, mouse, pencil, and many animals.³ Published in 2012, AlexNet was the first major deep convolutional network to leverage GPUs and is widely considered to be the spark that started the latest AI revolution.

The benchmark trains the AlexNet model using the [ImageNet](#) data set, a de-facto standard for deep learning training. The 143GB ImageNet data set contains 14,197,122 images from 21,841 distinct categories.

To reflect real AI development scenarios, we enabled distortion (image pre-processing steps). We also replicated the data by applying ten random data augmentation techniques to each JPEG image resulting in a 1.4TB data set and more than 141 million images. The 1.4TB data set was too large to fit in memory, forcing the system to repeatedly fetch data from the Isilon F800, ensuring that the benchmark stressed and measured entire system performance including compute, network, and storage.

To determine the maximum performance and scalability of the system while training the model, we ran the AlexNet benchmark multiple times, varying the number of GPUs and recording performance metrics of interest. Figure 6 shows the number of images processed per second while training AlexNet. Also shown are previously published results from a vendor with a solution that combines servers, NVIDIA Tesla V100 GPUs, and the vendor's custom all-flash storage solution.

³ <https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>

Figure 6. Dell EMC Ready Solutions for AI, Deep Learning with NVIDIA AlexNet Computational Throughput


Source: Enterprise Strategy Group

What the Numbers Mean

- With four GPUs, the Dell EMC Ready Solutions for Deep Learning processed almost 7,700 images/sec, taking a little more than five hours to complete, and processing 1.9 times more images/sec than Vendor A.
- With eight GPUs, the Dell EMC Ready Solutions for Deep Learning processed more than 14,250 images per second, cutting analysis time down to 2.77 hours, processing 2.9 times more images/sec than Vendor A.
- When the system was scaled to 16 GPUs, it processed more than 27,000 images per second, completing in just 86.5 minutes. Note: Vendor A did not publish results for the 16GPU benchmark variant.

AI algorithms benefit from parallel processing, and organizations can accelerate model training by using more processors in parallel. The ability to maintain data throughput and processing speed as the system scales to include more GPUs—the scaling efficiency—ensures that organizations maximize the return on their investment in bringing additional processing power to bear when training AI models. Table 1 shows the scaling efficiency of each solution while training AlexNet.

Table 1. Ready Solutions for AI, Deep Learning with NVIDIA Scaling Efficiency with AlexNet

GPUs	Dell EMC Ready Solutions for AI Deep Learning with NVIDIA	Vendor A
4	1.00	1.00
8	0.93	0.58
16	0.89	

Source: Enterprise Strategy Group

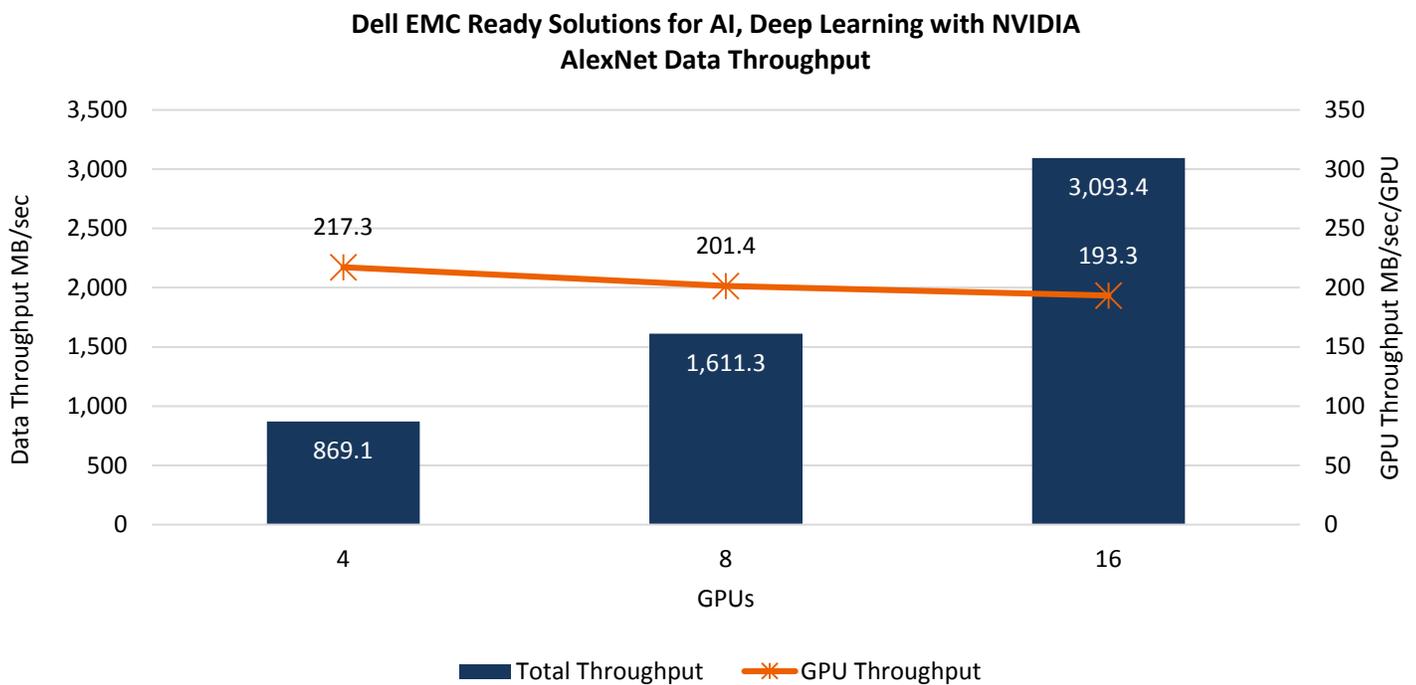
What the Numbers Mean

- Dell EMC Ready Solutions for AI, Deep Learning with NVIDIA scale efficiently, maintaining 93% of baseline performance (four GPUs) as the system was doubled to eight GPUs.

- Deep Learning with NVIDIA was almost as efficient as the system size was quadrupled, achieving 89% of its baseline performance when scaling from four to 16 GPUs.
- Vendor A’s solution was not nearly as efficient and was only able to achieve 58% of its baseline performance as the system size was doubled from four to eight GPUs.

Deep neural networks can have millions or even hundreds of millions of parameters (P). As a rule of thumb, ensuring a model’s ability to generalize (provide high accuracy predictions for any input) requires P^2 data points. Thus, organizations use multi-terabyte or even petabyte-sized data sets to train deep learning models, and the AI infrastructure needs to maximize and scale performance of the storage and data transport systems in addition to maximizing the raw computing power. Figure 7 shows the data throughput and throughput per GPU for AlexNet training using the Ready Solutions for AI, Deep Learning with NVIDIA system.

Figure 7. Dell EMC Ready Solutions for AI, Deep Learning with NVIDIA AlexNet Data Throughput



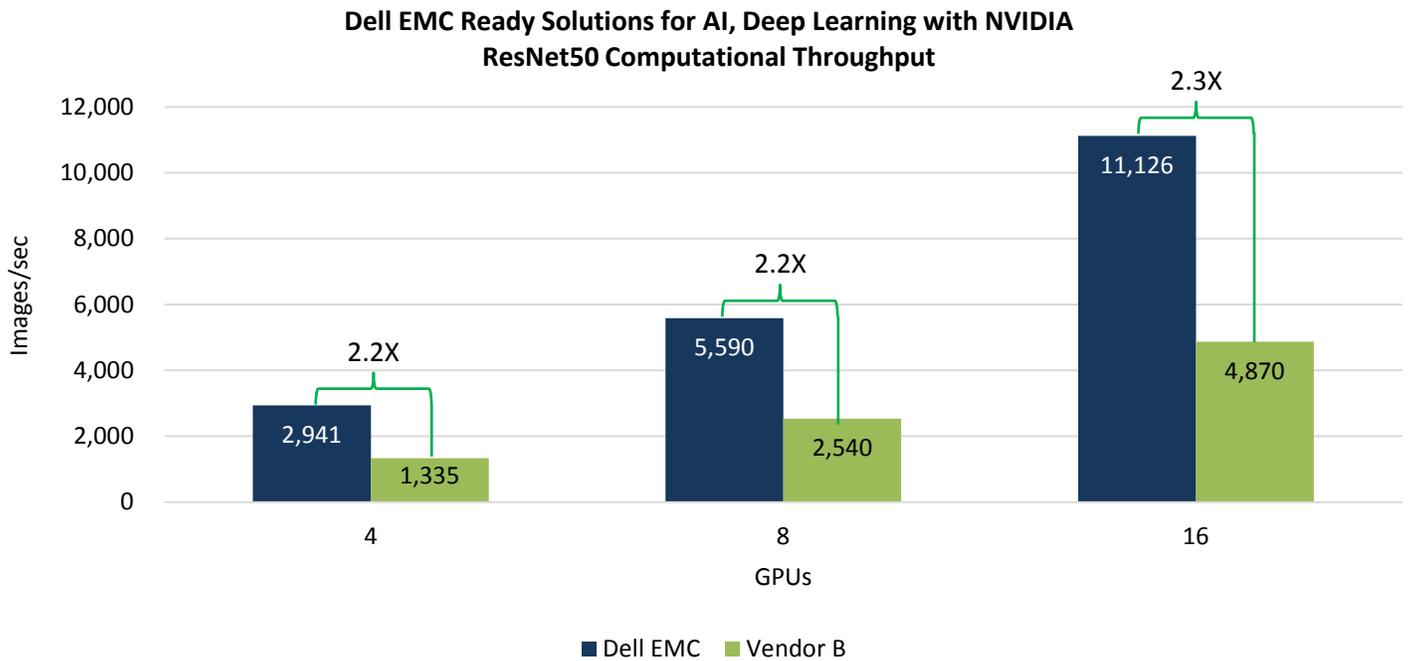
Source: Enterprise Strategy Group

What the Numbers Mean

- With four GPUs, Ready Solutions for AI, Deep Learning with NVIDIA transferred almost 870 MB/sec from the Isilon F800. As the solutions scaled to 16 GPUs, the solution transferred more than 3,000 MB/sec from storage.
- Throughout the tests, the GPUs averaged nearly 95% saturation. This high utilization shows that the Deep Learning with NVIDIA solution featuring Isilon is ideally architected to scale out and eliminate I/O bottlenecks for high bandwidth AlexNet training workloads.
- The solution transferred 217 MB/sec to each GPU in a four GPU solution. With up to 15 GB/sec of throughput in a single Isilon chassis, the Deep Learning with NVIDIA system can fully saturate an AlexNet workload up to 64 GPUs per Isilon chassis. Adding additional Isilon nodes linearly increases storage performance to support additional GPUs. Theoretically, with a maximum bandwidth of 540 GB/sec per cluster, the Isilon F800 can support 2,845 GPUs to process AI tasks similar to AlexNet. Note: The actual number of GPUs supported by Isilon will vary based on algorithm type, workload type, and data set size.

Next, we stressed Deep Learning with NVIDIA using ResNet50, a highly accurate image classifier published in 2015 by Microsoft research.⁴ The benchmark trains the ResNet50 model, which is much more computationally complex than AlexNet, using the same 1.4TB ten-times replicated ImageNet data set as was used for the AlexNet benchmark. Figure 8 shows the number of images processed per second while training ResNet50. Also shown are published results from a vendor with a solution that combines servers, NVIDIA Tesla V100 GPUs, and the vendor’s custom all-flash storage solution.

Figure 8. Ready Solutions for AI, Deep Learning with NVIDIA ResNet50 Computational Throughput



Source: Enterprise Strategy Group

What the Numbers Mean

- With four GPUs, Deep Learning with NVIDIA processed 2.2 times more images per second than Vendor B.
- Dell EMC’s performance advantage was maintained as the solution was scaled to use eight GPUs. Deep Learning with NVIDIA processed 2.2 times more images/sec than Vendor B.
- Dell EMC’s performance advantage was maintained as the solution was scaled to use 16 GPUs. Deep Learning with NVIDIA processed 2.3 times more images per second than Vendor B.

Table 2 shows the scaling efficiency of each solution while training ResNet50.

Table 2. Ready Solutions for AI, Deep Learning with NVIDIA Scaling Efficiency with ResNet50

GPUs	Dell EMC Ready Solutions for Machine Learning	Vendor B
4	1.00	1.00
8	0.95	0.95
16	0.95	0.91

Source: Enterprise Strategy Group

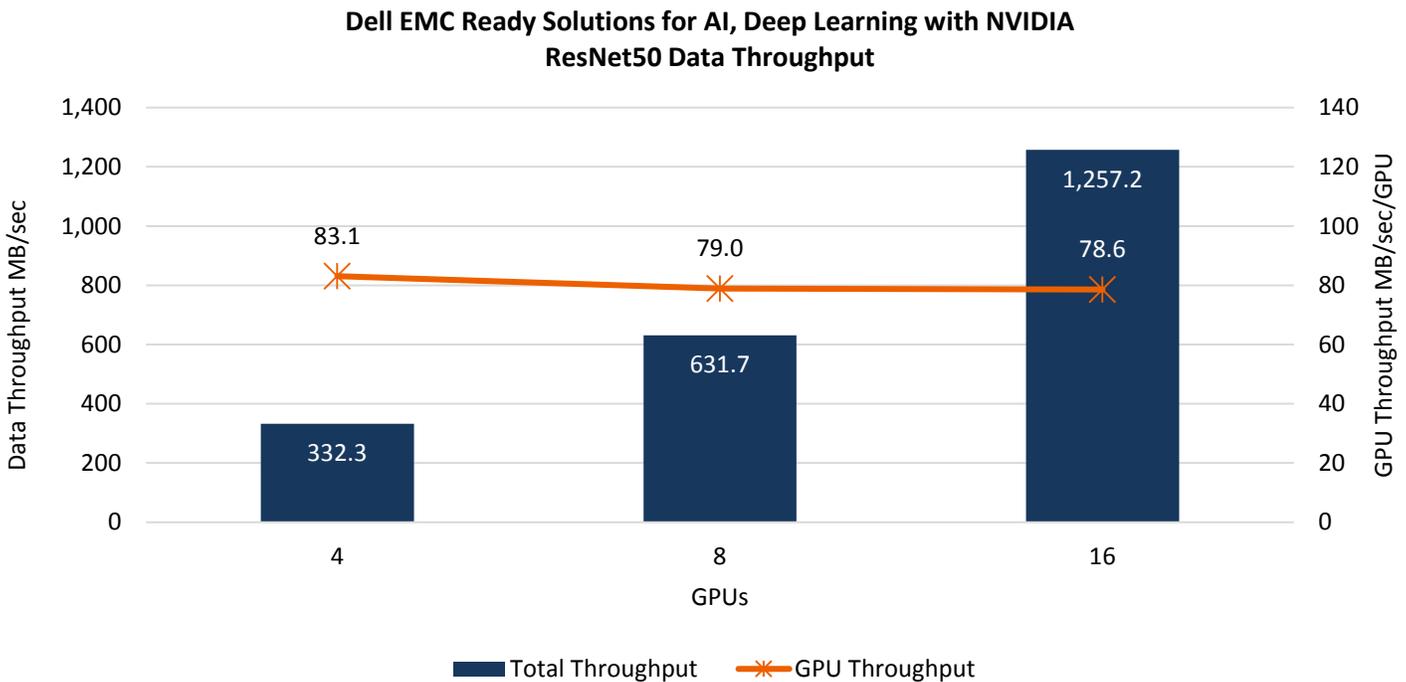
⁴ <https://arxiv.org/abs/1512.03385>

What the Numbers Mean

- Dell EMC Ready Solutions for AI, Deep Learning with NVIDIA scale efficiently, maintaining 95% of baseline performance as the system was doubled and quadrupled in size (eight and 16 GPUs).
- Vendor B’s solution demonstrated decreasing efficiency as the solution scaled—as it scaled from four to eight GPUs, it achieved 95% of baseline performance. However, as the solution scaled to 16 GPUs, it was only able to achieve 91% of baseline performance.

Figure 9 shows the system data throughput and throughput per GPU for ResNet50 training using Deep Learning with NVIDIA.

Figure 9. Dell EMC Ready Solutions for AI, Deep Learning with NVIDIA ResNet50 Data Throughput



Source: Enterprise Strategy Group

What the Numbers Mean

- With four GPUs, the Dell EMC Ready Solutions for AI, Deep Learning with NVIDIA transferred more than 330 MB/sec from the Isilon F800. As the solutions scaled to 16 GPUs, more than 1,250 MB/sec was transferred.
- Throughout the tests, the GPU's averaged nearly 95% saturation. This high utilization shows that the Deep Learning with NVIDIA solution featuring Isilon is ideally architected to scale out and eliminate I/O bottlenecks for ResNet training workloads.
- The solution transferred 83 MB/sec to each GPU in a four GPU solution and 78.6 MB/sec to each GPU in a 16GPU solution. With 15GB/sec throughput, the base Isilon F800 can saturate 180 GPUs. Adding additional Isilon modules increases storage bandwidth and concurrent connections to support additional GPUs. Theoretically, with a maximum bandwidth of 540 GB/sec, the Isilon F800 can support 6,500 GPUs to process AI tasks similar to ResNet50. Note: The actual number of GPUs supported by Isilon will vary based on algorithm type, workload type, and data set size.



Why This Matters

For AI, more complex models trained with larger data sets provide better results. With data sets in the tens of TBs to tens of PBs range and models with millions of parameters, high performance, high concurrency, and scale-out compute and storage become critical factors for organizations looking to obtain timely results from their AI efforts.

ESG validated that the 16GPU Dell EMC Ready Solutions for AI, Deep Learning with NVIDIA system was able to train the AlexNet model at 27,375 images/sec and the ResNet50 model at 11,126 images/sec. The scale-out Dell EMC solution featuring Isilon proved to be 2.2-2.9 times faster than systems from two other vendors. ESG also validated that Deep Learning with NVIDIA maintained processing speed as the system was scaled, achieving 89-95% of baseline performance as the number of GPUs was doubled and quadrupled. This ensures that organizations can maximize their return on investment as they scale out compute and storage to accelerate AI model development.

The Bigger Truth

Organizations perceive that AI is the next technology that will enable the faster delivery of better business outcomes. According to recent ESG research, 69% of respondents expect that ML and AI will deliver significant measurable outcomes in the near term, with 17% of respondents indicating that AI and ML were critical to their organization's strategy.

Lacking a standardized AI infrastructure stack, organizations can invest the time, effort, and money to select, acquire, integrate, configure, test, and validate their own custom stack. This complex process can take months, and the organization must juggle purchasing and support across many vendors. Public cloud solutions suffer from huge cost variability and the time and cost necessary to transfer and store terabytes to petabytes of data.

Dell EMC created the Ready Solutions for AI as standardized infrastructure stacks for machine learning and deep learning. These are validated and integrated hardware and software stack solutions, tuned and optimized to accelerate AI initiatives, shortening deployment time from months to weeks. Ready Solutions for AI simplify and accelerate data scientists' efforts, providing self-service workspaces where each data scientist can configure her own environment from a library of AI models and frameworks in just five clicks.

ESG validated that these solutions can accelerate AI model development. With PowerEdge C4140 servers accelerated with (16) NVIDIA GPUs and a chassis of Isilon F800 All-flash Scale-out NAS, Dell EMC Ready Solutions for AI trained the AlexNet model at 27,735 images/sec, and the more computationally complex ResNet50 model at 11,126 images per second. These results were from 2.2 to 2.9 times faster than results published by other vendors.

These integrated solutions for AI demonstrated scaling efficiency, keeping the GPUs pegged at 95% utilization while achieving 89-95% of baseline performance as the systems were scaled from four to eight to 16 GPUs, maximizing the return on investment as more GPUs are applied to solve more complex problems with larger and larger data sets. This high GPU utilization and linear scaling shows that the Deep Learning with NVIDIA solution featuring Isilon is ideally architected to scale out and eliminate I/O bottlenecks for AI training workloads.

ESG recommends that organizations investigate how Dell EMC Ready Solutions for AI can simplify and accelerate their AI journey.

All trademark names are property of their respective companies. Information contained in this publication has been obtained by sources The Enterprise Strategy Group (ESG) considers to be reliable but is not warranted by ESG. This publication may contain opinions of ESG, which are subject to change. This publication is copyrighted by The Enterprise Strategy Group, Inc. Any reproduction or redistribution of this publication, in whole or in part, whether in hard-copy format, electronically, or otherwise to persons not authorized to receive it, without the express consent of The Enterprise Strategy Group, Inc., is in violation of U.S. copyright law and will be subject to an action for civil damages and, if applicable, criminal prosecution. Should you have any questions, please contact ESG Client Relations at 508.482.0188.

The goal of ESG Validation reports is to educate IT professionals about information technology solutions for companies of all types and sizes. ESG Validation reports are not meant to replace the evaluation process that should be conducted before making purchasing decisions, but rather to provide insight into these emerging technologies. Our objectives are to explore some of the more valuable features and functions of IT solutions, show how they can be used to solve real customer problems, and identify any areas needing improvement. The ESG Validation Team's expert third-party perspective is based on our own hands-on testing as well as on interviews with customers who use these products in production environments.