Technical Validation

# Big Data Analytics Applications Using HPE Ezmeral Unified Analytics with NVIDIA RAPIDS Accelerator

## Optimizing Performance and Management of Apache Spark Containerized Workloads

By Alex Arcilla, Senior Validation Analyst

December 2021

## Introduction

This ESG Technical Validation documents our evaluation on HPE Unified Analytics. We focus on how the combination of HPE Ezmeral Runtime Enterprise and NVIDIA RAPIDS Accelerator can accelerate Apache Spark AI/ML workloads running on NVIDIA GPU-enabled containers. We also reviewed how organizations can use HPE Unified Analytics to simplify management and orchestration of multiple containerized applications simultaneously.

### Background

According to ESG research, 34% of survey respondents cited providing more predictive insights into future scenarios or outcomes as among their most important objectives when investing in AI/ML initiatives, while another 34% cited developing AI-based products/services to capture new revenue opportunities (see Figure 1).[1] However, challenges remain in implementing an infrastructure that will achieve these objectives. The same ESG research revealed that resource sharing is the most-cited part of the infrastructure stack that respondents believe is one of the weakest links in delivering an effective AI environment.[2]

**Figure 1.  Most Important Objectives to Accomplish with AI/ML Investments**



What are the most important objectives you expect to accomplish from your organization's investments in the area of AI/ML? (Percent of respondents, N=325)

| | |
|---|---|
| Providing more predictive insights into future scenarios or outcomes | 34% |
| Developing AI-based products/services to capture new revenue opportunities | 34% |
| Improving customer experience | 34% |
| Reducing risk around business decisions and strategy | 33% |
| Enabling faster tactical response to shifting customer requirements | 33% |

*Source: Enterprise Strategy Group*

Current AI/ML initiatives have been implemented differently across different business units, either on-premises or in the public cloud, using a wide variety of on-premises infrastructure, public cloud infrastructure and services, open-source tools, container engines, and libraries. This has resulted in multiple and disjointed AI/ML environments deployed across the organization. However, multiple data silos and compute resources remain separate, limiting how they can be fully utilized across the organization for maximum business advantage.

What if there was a solution that can maximize the usage of AI/ML compute and storage resources across the organization, whether those resources are located on-premises or in the public cloud, without sacrificing the performance required to gain insight and meet business needs?

---

[1] Source: ESG Survey Results, *Supporting AI/ML Initiatives with a Modern Infrastructure Stack*, May 2021.
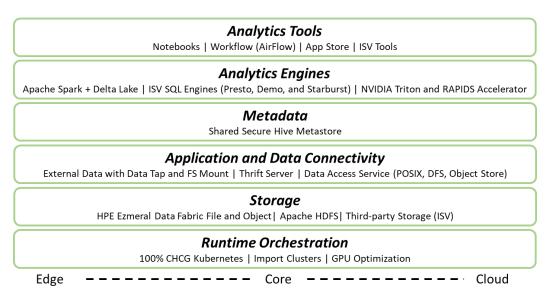
[2] Ibid.

## HPE Ezmeral Unified Analytics with NVIDIA RAPIDS Accelerator

HPE Ezmeral Unified Analytics, available in January 2022, has been designed to deliver a unified data lake that connects distributed data stores, public cloud storage, files, and object stores, regardless of where they reside. The foundation of HPE Ezmeral Unified Analytics is HPE Ezmeral Runtime Enterprise, which delivers containerized orchestration of both cloud-native and non-cloud applications with a persistent data store. HPE Ezmeral Runtime Enterprise is designed to help organizations simplify the development, orchestration, and scaling of Kubernetes-based clusters deployed across hybrid environments.

The platform can help data scientists, data engineers, and data analysts to support multiple enterprise-grade containerized AI/ML applications via the following capabilities (see Figure 2):

**Figure 2. HPE Unified Analytics**



> **Analytics Tools**
> Notebooks | Workflow (AirFlow) | App Store | ISV Tools
>
> **Analytics Engines**
> Apache Spark + Delta Lake | ISV SQL Engines (Presto, Demo, and Starburst) | NVIDIA Triton and RAPIDS Accelerator
>
> **Metadata**
> Shared Secure Hive Metastore
>
> **Application and Data Connectivity**
> External Data with Data Tap and FS Mount | Thrift Server | Data Access Service (POSIX, DFS, Object Store)
>
> **Storage**
> HPE Ezmeral Data Fabric File and Object| Apache HDFS| Third-party Storage (ISV)
>
> **Runtime Orchestration**
> 100% CHCG Kubernetes | Import Clusters | GPU Optimization
>
> Edge – – – – – – – – – – – – – Core – – – – – – – – – – – – – Cloud

*Source: Enterprise Strategy Group*

- **Multi-tenancy**: With the platform's control plane, organizations can spin up and orchestrate multiple standards-based Kubernetes clusters using physical or virtual compute (CPU or GPU) and storage resources located on-premises or in the public cloud. Organizations can also automate orchestration of workflows, allowing many experiments to be scheduled and run "hands free." Thus, data science and development teams can share the same infrastructure stack with minimal impact on system performance.

- **Unified data lake**: The single global namespace provided by HPE Ezmeral Data Fabric provides a consolidated view into files that can reside in separate clusters across on-premises, cloud, and edge environments, thus enabling end-users and their applications to access remote data as if it is locally stored.

- **Open platform/ecosystem**: HPE Ezmeral Unified Analytics is built on an open source foundation that enables data analytics teams to deploy Kubernetes-based workloads into any AI/ML environment, comprising a wide range of open source tools, container engines, and libraries. NVIDIA Triton, RAPIDS Accelerator, and Apache Spark are included in the platform's ISV marketplace (see Figure 10 in the Appendix for details).

- **Security**: Organizations can integrate HPE Ezmeral Runtime Enterprise with enterprise LDAP and AD authentication systems to secure access to jobs, data, or clusters based on groups and/or roles. Access to the multiple applications

supported by the platform can be restricted according to business needs. Risk and compliance tenant structure allows for granular security structures.

To simplify platform use, multiple versions of Kubernetes are supported, either to accommodate differing levels of expertise throughout the organization or to maintain current procedures and toolsets without disrupting business operations. End-users can use either API commands or a graphical web UI with this self-service platform, removing the complexity typically faced when using Kubernetes.

For ML initiatives, organizations have typically relied on GPU-based infrastructure to accelerate processing time, given the large data sets required to obtain a useful ML model. Organizations that employ Kubernetes clusters based on NVIDIA's A100 Tensor Core GPUs can leverage NVIDIA RAPIDS Accelerator for Apache Spark. While Apache Spark is designed to increase data processing performance via in-memory processing, NVIDIA RAPIDS Accelerator for Apache Spark can maximize GPU performance. RAPIDS is a suite of open-source software libraries and APIs for executing data science pipelines entirely on GPUs. The software has been designed to accelerate data preparation, feature extraction, training, and inference with model development and training.

## ESG Technical Validation

ESG evaluated how HPE Ezmeral Runtime Enterprise and NVIDIA RAPID's Accelerator can accelerate the performance of Apache Spark workloads. We also observed the simplicity of Kubernetes cluster creation and the ability to access data regardless of location. Testing was performed via remote sessions.
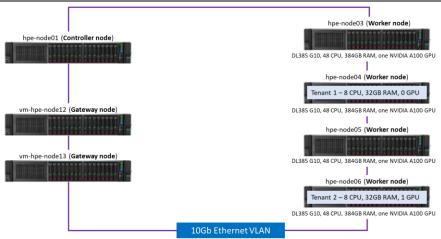
### Performance

Apache Spark-based ML workloads require fast processing and response times to accommodate the very large data volumes required to create and train ML models. The combination of HPE Ezmeral Unified Analytics, NVIDIA RAPIDS Accelerator, and NVIDIA A100 GPU-based clusters can accelerate Apache Spark workloads, which can benefit applications such as fraud detection or financial transaction processing.

### ESG Testing

The following test bed was used (shown in Figure 3). Six HPE ProLiant servers (one master node, one controller node, and four worker nodes) were used. Each worker node was equipped with one AMD EPIC 7352 24-core CPU and one NVIDIA A100 40G GPU. The master and controller node were equipped with two AMD EPYC 7702 64-core CPUs. All nodes were interconnected by a 10Gb Ethernet VLAN configured on a Cisco Nexus Model 9000 switch.

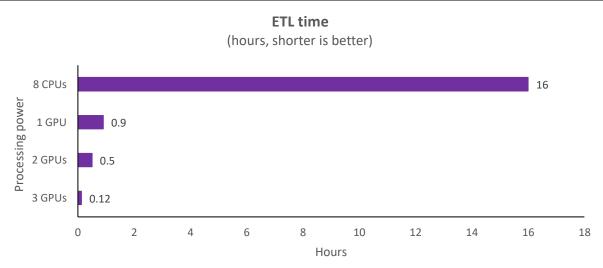**Figure 3. Test Bed for Apache Spark Workload Performance Testing**



*Source: Enterprise Strategy Group*

Clusters were created with HPE Ezmeral Runtime Enterprise and Kubernetes v1.20. Container size for CPU-based workloads was limited to eight CPUs and 32GB RAM and one container with 48 CPUs, 384 GB of RAM, and one NVIDIA A100 GPU. Two separate namespaces were configured for the CPU and GPU-based clusters, respectively. All jobs were automatically submitted via Apache Spark Submit. Both extract-transform-load (ETL) and query response times were measured using a 5TB data set.

ESG first reviewed the test results that measured ETL processing time, covering the creation of data sets and tables, as well as conversion of file format to Apache Parquet. ETL times were recorded when processing data sets on one container (equivalent to eight CPUs) were measured against another container with GPUs scaling from one, two, and three NVIDIA A100 GPUs respectively. The NVIDIA RAPIDS Accelerator ran when using GPU-based containers. Results are shown in Figure 4.

**Figure 4. ETL Times with Apache Spark and NVIDIA RAPIDS Accelerator on HPE Ezmeral Runtime**



*Source: Enterprise Strategy Group*

ESG then reviewed the results of query processing time. We selected four jobs—Queries #5, #16, #21, and #22—to illustrate performance. Each job ran on four Kubernetes clusters with container sizes equivalent to eight CPUs and one,

two, and three GPUs, respectively. Again, the NVIDIA RAPIDS Accelerator ran when using GPU-based clusters. We compared query response times, shown in Figure 5 and Table 1.

**Figure 5. Query Response Times with Apache Spark and NVIDIA RAPIDS Accelerator on HPE Ezmeral Runtime**



*Source: Enterprise Strategy Group*

**Table 1. Query Processing Times Using Apache Spark and NVIDIA RAPIDS Accelerator**

|  | Query #5 Query response time (min) | Query #16 Query Response Time (min) | Query #21 Query Response Time (min) | Query #22 Query Response Time (min) |
|---|---|---|---|---|
| 8 CPUs | 32.13 | 2.31 | 1.1 | 0.5 |
| 1 GPU | 12.15 | 1.2 | 0.5 | 0.35 |
| 2 GPUs | 9.13 | 0.56 | 0.25 | 0.12 |
| 3 GPUs | 5.16 | 0.18 | 0.11 | 0.09 |

*Source: Enterprise Strategy Group*

## What the Numbers Mean

- ESG confirmed that HPE Ezmeral Runtime Enterprise with NVIDIA RAPIDS Accelerator, and NVIDIA A100 GPU-enabled Kubernetes clusters can improve performance of Apache Spark ML workloads by up to 29x compared to running similar workloads on an eight CPU-enabled cluster. (Results will vary in real-world production environments.)

- ESG also found that when leveraging the NVIDIA RAPIDS Accelerator, query processing times on a single GPU-based cluster decreased between 90% and 97%. Response times decreased as more GPUs were added to the cluster.

- Optimizing Apache Spark workload performance translates into minimizing both ETL and query response times. This is critical to helping organizations quickly process large data sets to deliver business insights quickly.

## ⓘ Why This Matters

While organizations embracing containerized applications find it important to simplify their deployment and management, application performance must be maintained. This is especially critical for those wanting to use Kubernetes clusters to support big data and AI/ML applications and workloads, which require high-performance and large-scale data processing.

ESG confirmed that processing of Apache Spark workloads was accelerated with the addition of HPE Ezmeral Runtime Enterprise, NVIDIA RAPIDS Accelerator and NVIDIA A100 GPUs. Customers can experience even greater acceleration when NVIDIA A100 GPU capacity is increased. We saw how the combination of HPE Ezmeral, NVIDIA, and Apache Spark allows enterprise customers to get insights faster for more data-driven business decisions.

## Simple to Use

With HPE Ezmeral Runtime Enterprise, organizations can create and manage multiple Kubernetes clusters and containerized applications within a single web UI, without the need to purchase or configure the underlying infrastructure and supporting toolsets.

### ESG Testing

ESG began by observing the creation of a Kubernetes cluster using a wizard-driven process (see Figure 6). For this cluster, we used Amazon Web Services (AWS) Elastics Compute Cloud (EC2) instances as hosts. We navigated to the HPE Ezmeral Runtime Enterprise web UI; selected the virtual hosts for the cluster; assigned AWS EC2 instances as master and worker nodes; selected the Kubernetes version to run, IP network address ranges, LDAP and AD settings; and added other services, such as Apache Spark Operator.
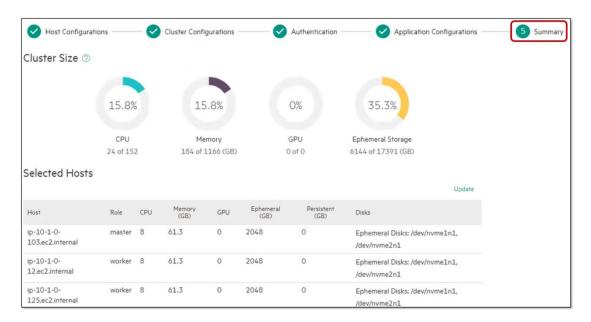
**Figure 6. Creating a Kubernetes Cluster with HPE Ezmeral Container Platform Web UI**



Once we completed this process, we viewed the details of the Kubernetes Cluster in a Summary screen (see Figure 7). The summary detailed the amount of compute and storage capacity available to any application running on this cluster and noted the specific physical and virtual hosts consumed by this cluster.

**Figure 7. Summary Details of Kubernetes Cluster**



ESG observed that this process completed in minutes, resulting in significant savings in both time and effort. Creating a Kubernetes cluster on premises involves many steps. Hardware needs to be purchased and configured with the proper compute and storage capacity to support gateway, controller, master, and worker nodes. Networking and security requirements need to be addressed. All open source software, especially Kubernetes and supporting tools, need to be downloaded from code repositories. Furthermore, additional scripts need to be written to customize the software to satisfy specific business needs. The net impact of all these steps is that time, effort, and cost are wasted, and time to value is long.

Deploying containerized applications on Kubernetes clusters may be simpler when working with public cloud providers. However, considering the cost of virtual compute and storage resources, services required for containerization and cluster management, this may not be a cost-effective alternative for many organizations. Also, time and effort can be excessive as organizations need to coordinate tasks between multiple interfaces, regardless of the public cloud provider.

> **i Why This Matters**
>
> ESG validated that HPE Ezmeral Runtime Enterprise reduces the amount of time to operationalize Kubernetes clusters with wizards and automated prechecks to deliver business outcomes faster.
>
> We observed the ease and simplicity of creating a cluster, removing the multiple tasks and complexity of setting up the infrastructure on premises or using multiple public cloud-based services concurrently.

## Access Data Regardless of Data Location

Application performance partly depends upon how quickly data can be accessed. However, replicating and transmitting data to local storage may not always be an option. For example, organizations want to analyze data from IoT deployments in real time. That cannot happen when data is being gathered at the network edge while the analytics application is located within a data center. Applications placed into containers by HPE Ezmeral Runtime Enterprise can ingest any data type or source, regardless of the data's physical location.

### ESG Testing

ESG reviewed the deployment of a Kubernetes cluster using a Jupyter notebook and HPE Ezmeral Runtime Enterprise. As shown in Figure 8, this application was designed to predict power output from wind turbines in Australia. The application used Apache Spark for data processing and accessed data via a secure connection to HPE Ezmeral Data Fabric, enabling end-users to ingest data from remote IoT or cloud deployments.

We observed that a few lines of declarative script that directed path names to data sources were all that was required for this application to access the IoT data. The data fabric was also deployed as a Kubernetes cluster, allowing it to transparently scale as data is ingested.
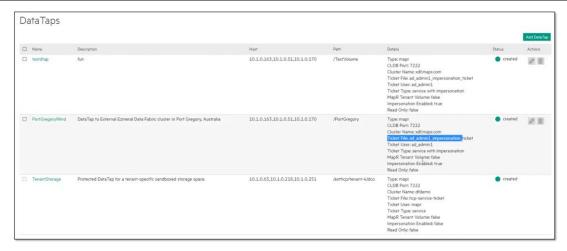
**Figure 8. Jupyter Notebook Application Accessing IoT Data**



Creating secure connections to IoT data was also simple, as shown in Figure 9. Using another wizard-driven process, we observed the simplicity of creating secure connections to any data source by inputting key information such as the IP addresses of the hosts collecting the data.

**Figure 9. Details of Embedded Connections to IoT Data**



We ran this Jupyter application without replicating the data to a local data store or impacting application performance.

> ### ℹ️ Why This Matters
>
> Businesses need to run analytics in-place where the data is located. This prevents hurdles that slow down insights due to the need to move or manipulate the data from an edge or cloud environment.
>
> ESG verified that HPE Ezmeral Runtime Enterprise uses APIs and POSIX to allow developers and data scientists to use script language and path names to access data wherever it is located. This saves the time and effort of copying data to secondary systems before analysis can begin.

## The Bigger Truth

As more organizations employ AI/ML applications to gain business advantage, they are faced with a wide variety of infrastructure and tools for implementing an AI/ML environment. However, the options available for building these environments are not standardized. Multiple business units end up implementing their own AI/ML infrastructure stacks, leading to disjointed compute and data resources that are not fully leveraged. Orchestrating and managing these diverse resources and the underlying infrastructure leads to operational inefficiency and decreased time to insight and value.

HPE Ezmeral Unified Analytics is an open, hybrid platform that unifies analytics across disjointed data stores to provide data teams with a single data repository, simplifying data analytics. HPE Ezmeral Runtime Enterprise, built into HPE Ezmeral Unified Analytics, eliminates steep learning curves, as well as the time and resources required to create a Kubernetes infrastructure across hybrid deployments.

ESG validated that HPE Ezmeral Runtime Enterprise, the foundation of HPE Unified Analytics, can:

- Accelerate performance of Apache-Spark workloads when using the combination of HPE Ezmeral Runtime Enterprise, NVIDIA A100 GPU-enabled Kubernetes clusters and NVIDIA RAPIDS Accelerator.

- Simplify the use of Kubernetes via wizard-enabled processes to remove the complexity of creating and managing Kubernetes clusters.

- Enable data access regardless of data location, thus not sacrificing overall application performance.

Potential issues that ESG believes organizations should examine when considering the use of HPE Ezmeral Unified Analytics include application size and the number and type of applications to be containerized. Small applications are not the best candidates for containerization; better candidates skew toward web-scale applications. Also, the platform may be better suited for organizations planning to modernize multiple applications (both re-factoring existing ones and net-new).

If your organization seeks to simplify how containers are deployed and managed while maximizing performance of workloads running on NVIDIA A100 GPU workloads, ESG recommends taking a closer look at how HPE Ezmeral Unified Analytics, NVIDIA RAPIDS Accelerator and NVIDIA A100 GPUs.

## Appendix

### Figure 10: ISV Marketplace

**ESG**

**Enterprise Strategy Group** is an IT analyst, research, validation, and strategy firm that provides market intelligence and actionable insight to the global IT community.