**intel.**

# DAOS and Google Cloud Provide a Balanced and Performant HPC Scale-out Computational Environment

## Distributed Asynchronous Object Store (DAOS) provides extreme storage performance on Google Cloud.

The appeal of high performance computing (HPC) in the cloud is undeniable: spin-up an instance, run on the latest hardware, and get your results. The "supercomputer" is always there, ready to run 24/7, with no lengthy review process to get time or hardware to purchase. Even better, Google Cloud has automated the experience with a toolkit that gives users the flexibility to create and tear down entire HPC and AI clusters within minutes.

> With DAOS on Google Cloud, users can easily and quickly provision storage clusters that can scale to similar performance levels as similar on-premises hardware, but then also be able to dynamically grow and shrink those clusters as needed.
>
> — Andrey Kudryavtsev, Intel

The Google Cloud Platform (GCP) team noted at the Supercomputing 2022 (SC'22) conference that adding DAOS provides scalable performance as both IOR and IO500 distributed benchmark tests were able to achieve balanced and strong performance. Performance results of 96 GiB/s read, 60 GiB/s write, and IO latencies as low as 0.28ms random read and 0.36ms per random write were demonstrated in a distributed Google Cloud environment. Andrey Kudryavtsev, DAOS product manager at Intel, explained, "With DAOS on Google Cloud, users can easily and

## Table of Contents

| | |
|---|---|
| **Fast** | **High-bandwidth**, **high IOPs**, and **low-latency**, utilizing byte-addressable media to store metadata and small I/O, and local attached NVMe for bulk I/Os. |
| **Reliable** | Use **software managed redundancy** configurable per object to protect your data just when you need to, and not when you don't. |
| **Cost Effective** | **Cheaper per-BW and per-IOP** than comparable solutions. Combine with Google Cloud Storage for a cost-effective solution to meet your bandwidth and capacity needs. |
| **Next Gen** | Designed for **memory addressable storage** and supports access via **POSIX or a custom library that removes POSiX overheads** |

Figure 1. DAOS: extreme performance storage system on Google Cloud.

**Slurm** 1st
**Slurm on Google**
Auto-scaling, hybrid features built for Google

**DDN**
**DDN EXAScaler** 1st
DDN EXAScaler released on Google, first EXAScaler in cloud

**Dell EMC**
**Dell PowerScale**
storage solution launched on Google

**HPC VM Image**
Best practice CentOS 7 VM image with HPC & MPI tunings

**IBM LSF**
Improvements to the IBM LSF Google Cloud Resource Connector

**Filestore High Scale**
Scalable managed NFS service, now supporting up to 100TiB and 26GB/s

**HPC Toolkit**
Simple best practices toolkit to deploy HPC environments on Google

**Transfer Service On-Prem**
Simple API-driven data transfer tool for GCS and file storage systems

**C3 VMs**
Intel Sapphire Rapids VMs for compute-intensive workloads

**2019** | **2020** | **2021** | **2022**

**T4** 1st
**NVIDIA GPUs**
NVIDIA T4 GPUs on N1 VMs; up to four per VM

**C2** 1st
**intel VMs**
Intel Cascade Lake VMs. 3.8GHz turbo clock, 40% higher performance vs N1

**A2** 1st
**NVIDIA VMs**
16 NVIDIA A100 GPUs per VM, NVLink enabled

**Placement Policies**
Compact / Spread network placement. Compact for ~20 VMs per placement group.

**100Gbps, gVNIC & 9K MTU**
NIC Driver & 9K MTU and 100Gbps enables high throughput, low latency networking between VMs

**Bulk API**
Create up to 1,000 VMs per API call with regional capacity finding

**AMD C2D VMs**
AMD EPYC Milan VMs. High clock speed, large memory and high vCPU count VMs.

**Batch**
Simple cloud-native batch scheduler for running workloads in a managed, zero-ops service

**T2A VMs**
ARM-based VMs optimized for scale-out workloads

—●— Generally Available
●····· In Preview
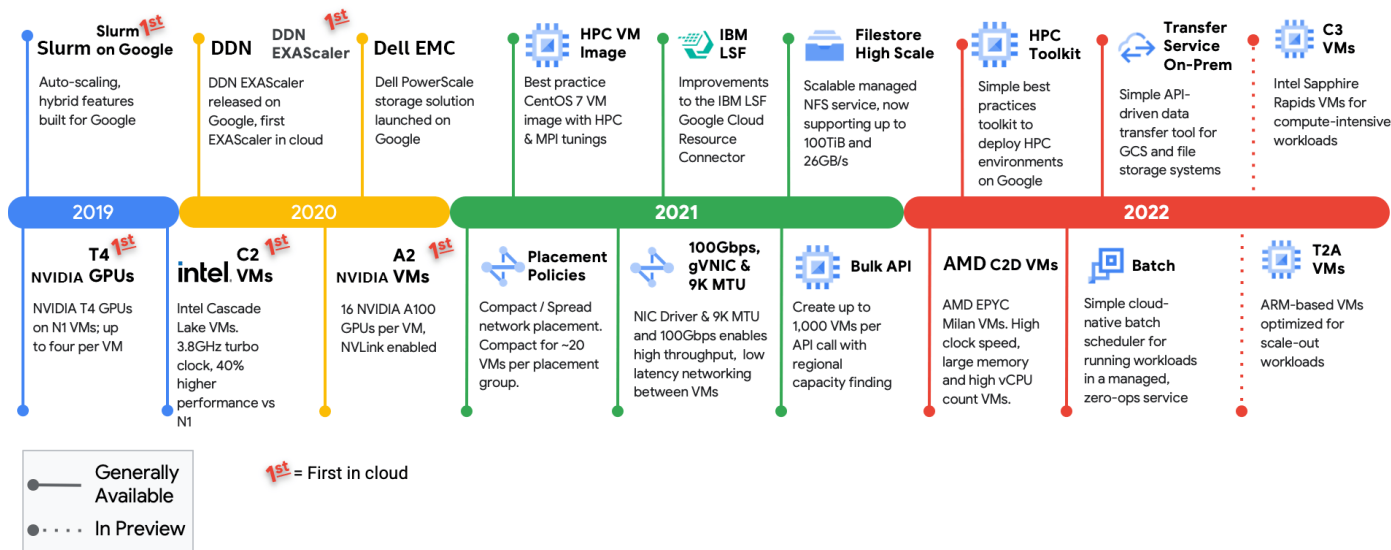
1st = First in cloud

Figure 2. Google has been investing in HPC infrastructure and software for years.

quickly provision storage clusters that can scale to similar performance levels as similar on-premises hardware, but then also be able to dynamically grow and shrink those clusters as needed."

*As part of our collaboration with Intel, the DAOS software is now part of the Google Cloud HPC ToolKit. Its inclusion reflects part of a greater opportunity for Google Cloud to address compute, communications, and even more storage intensive HPC workloads.*

— Carlos Boneti, Google Cloud

## The Realization of a Multi-Year Effort

Since 2019, Google has made a significant investment in HPC. The HPC Toolkit, Transfer Service, and 4th Gen Intel Xeon® processor-based virtual machines (VMs) (shown in the upper right of Figure 2) make it easy to migrate HPC applications to Google Cloud Platform (GCP) and use DAOS with the latest computational and storage technology. A VM represents a single node or client in the HPC cluster.) The new Google Cloud C3 VMs will be the first VM in the public cloud to run 4th Gen Intel Xeon processors.[1] For a sneak preview of DAOS performance on the latest 4th Gen Intel Xeon processors, see "Performance Evolution of DAOS Servers".

Carlos Boneti, (HPC Software Engineer, Google Cloud, highlights Google's investment in HPC, "Google has been building HPC solutions for years in partnership with many organizations. Our partners build products quickly, but often produce separate solutions that are focused on individual HPC challenges. The Cloud HPC Toolkit provides an ecosystem where multiple components or modules work together in a single blueprint that can be configured and instantiated in an automatic fashion, thus saving user time and errors. As part of our collaboration with Intel, the DAOS software is now part of the Google Cloud HPC ToolKit. Its inclusion reflects part of a greater opportunity for Google

Cloud to address compute, communications, and even more storage intensive HPC workloads."

The emphasis on storage reflects a recurring concern by the HPC community about data in the cloud, a challenge exacerbated by the recent popularity of data-intensive machine learning (ML) and high-performance data analytics (HPDA) workloads. For HPC-in-the-cloud, Kudryavtsev notes, "We are seeing great advantages for HPC customers in expanding to the cloud and fully leveraging its capabilities for HPC workloads. Google Cloud put significant effort into HPC Toolkit development to simplify this process for users who previously had to deal with complex cloud infrastructure manually. I am glad to see DAOS being a part of that integration. It makes perfect sense because of DAOS's storage capabilities focused on high performance, low cost per I/O, support for various fabrics including Ethernet (TCP and RoCE) and various cloud technologies (such an integration with containers). In less than two minutes, DAOS can be fully provisioned in GCP with the HPC Toolkit, serving the needs of the high-performance storage tier and co-existing with Google Cloud Storage for the best ROI for users."

*We are seeing great advantages for HPC customers in expanding to the cloud and fully leveraging its capabilities for HPC workloads. Google Cloud put significant effort into HPC Toolkit development to simplify this process for users who previously had to deal with complex cloud infrastructure manually. In less than 2 minutes, DAOS can be fully provisioned in GCP with the HPC Toolkit, serving the needs of the high-performance storage tier and co-existing with Google Cloud Storage for the best ROI for users.*

— Andrey Kudryavtsev, Intel

## Intel DAOS' World-Record Setting Scalable Storage Performance Moves to HPC-in-the-Cloud

Google Cloud has the hardware capability to speed the most computationally intensive HPC workloads with fast processors and access to GPU and TPU accelerators. It also has the storage capacity to service even the largest data sets, but HPC is well-known for finding the weakest link in any system or cloud instance. HPC workloads require a balanced computing environment where storage and communication performance have to be sufficient to service the needs of the computational nodes. Otherwise, bottlenecks can slow time-to-solution and increase the cost of cloud computing.

DAOS was recently recognized with the top bandwidth award at Supercomputing 2022. Intel DAOS was previously recognized as delivering world record setting storage performance.[2,3] These performance high water marks explain why DAOS is gaining momentum in the HPC community. In a flagship demonstration in support of extreme scale out workloads, the Intel-based Aurora exascale supercomputer will use DAOS as its primary storage system.[4,5] Use on Aurora ensures that the open-source DAOS software stack will be stress and reliability tested under the extreme workloads of a leadership-class exascale HPC supercomputer.

## A User-Space Solution

The open source DAOS storage software stack achieves world record setting performance by providing a storage solution that runs in user space.

This innovative architecture means that key HPC tools like HDF5, MPI-IO, TensorFlow, and more can call DAOS directly via a user space library (libdfs). Critically, DAOS also supports traditional HPC workloads that access their data through a file system mount point. Through use of an interception library, DAOS can support the traditional mount point yet eliminate the overhead of a kernel system calls and associated overhead to manage all the internal kernel filesystem and block caches. This makes DAOS an ideal solution for both born-in-the-cloud and traditional applications and HPC workloads in GCP.

Internally, the DAOS software is designed to eliminate many of the metadata and locking issues of traditional POSIX-based filesystems (Figure 3). These issues force serialization

and limit scalability and block performance. The result is high-bandwidth performance (even exceeding that of Lustre on some workloads) combined with ultra-low-latency storage access plus tremendous scalability.[8]

The SC'22 demonstration by the Google team highlighted many of the benefits of DAOS including its use of a key-value architecture, which avoids many POSIX limitations, and which differentiates DAOS from other storage solutions with its low-latency, built-in data protections, end-to-end data integrity, and overall increased single client performance. For these reasons, the Google team recommends that DAOS should be used for any workload where small file, small IO, and/or many metadata operations per second (MDop/s) performance is critical.

Kudryavtsev explains that this key-value design and other internal features of the DAOS architecture allow even the most extreme-scale Google Cloud installations: "The performance advantages of DAOS extend to the Google Cloud distributed environment as every node has the metadata associated with the data it owns (Figure 4). This makes DAOS fully distributed, with no single point metadata bottleneck and gives the node full control per dataset over how the metadata is updated. DAOS will most likely not be the limiting factor in scaling your HPC workload."
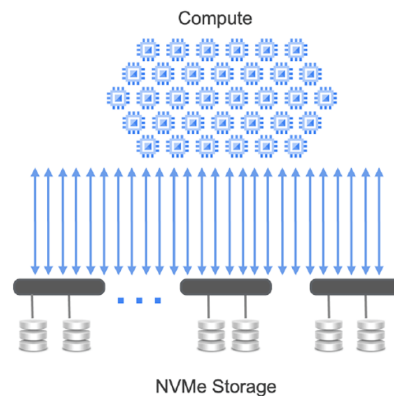


Figure 4. The DAOS architecture provided direct access to both data and metadata to deliver both high-performance and scalability.
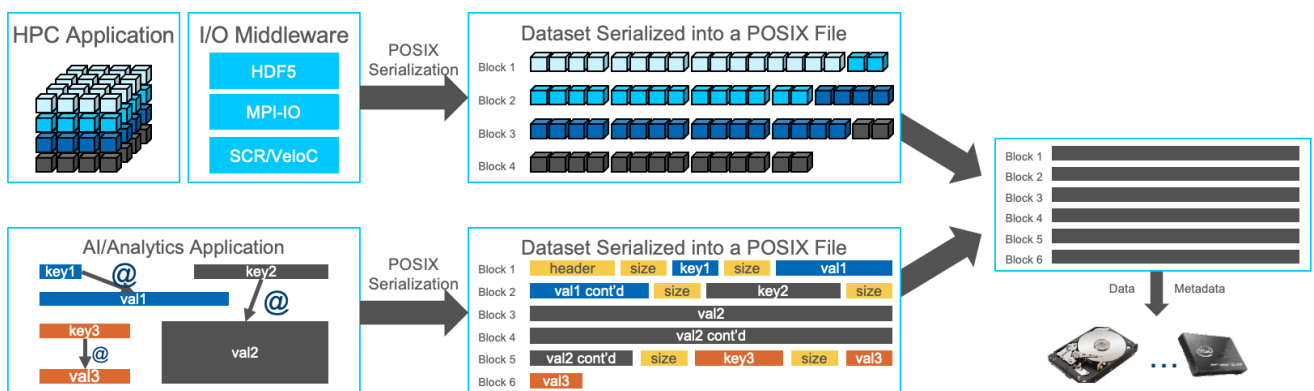


Figure 3. The problem with POSIX and blocks.

## Use Cases Highlight the Benefits of DAOS and Cloud Storage

The HPC software ecosystem on the Google Cloud is quite extensive as shown in Figure 5, which should support most user needs.

Data ingestion into Google Cloud is managed through the GCP Storage Transfer Service. A typical use case migrating on-premises data is illustrated below.

Barak Epstein, product manager for HPC storage at Google, explains, "Users can leverage Google Cloud Storage (GCS) for data ingestion from sources across the globe and long-term retention of that data at low cost. They can then transfer datasets of interest to DAOS for high performance analysis, drastically reducing the execution time (and hence the execution cost) of their high-performance apps. This cost- and performance-blended storage model helps users take advantage of cloud innovation while reducing the runtime and overall cost of both HPC and AI/ML applications."

**Following are typical HPC/AI/HPDA use cases:**

- **Traditional HPC:** Applications can use their HPC storage middleware of choice (e.g., HDF5, MPI-IO, POSIX) to accelerate I/O to DAOS. To lower cost, users store input datasets along with analysis results in GCS for long term retention.

- **AI:** Data-driven model development can require high I/O operations per second (IOP/s) and/or MDop/s. DAOS has demonstrated its capability to deliver both in the Google Cloud environment as well as on-premises. The high IOP/s, MDop/s, and low latency provided by DAOS is ideal to support AI/ML workloads with high I/O demands. One example is computer vision, a data-intensive workload with very high throughput requirements. Epstein notes, "DAOS allows the porting of AI/ML workloads to cloud by providing I/O performance not available elsewhere."
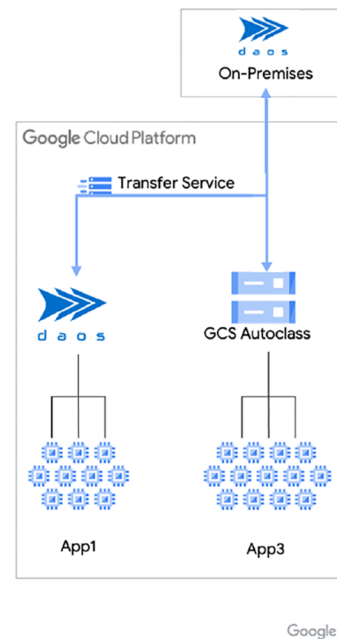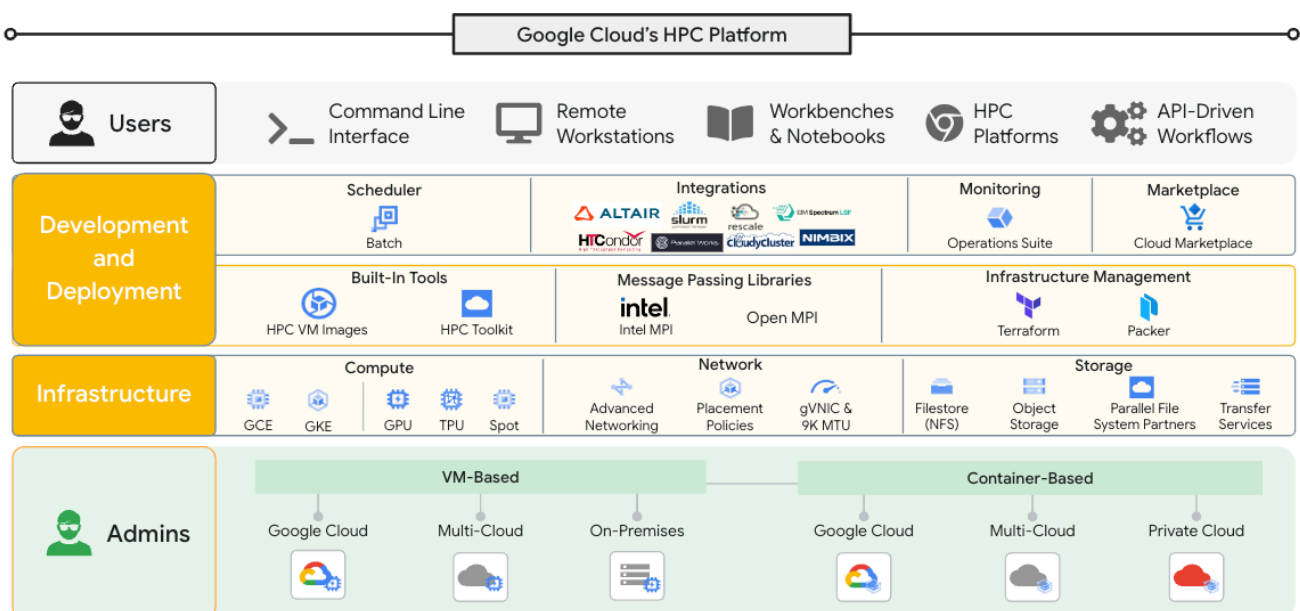


Figure 6. Building a hybrid storage layer with DAOS and GCS.

- HPDA: Many users are looking to Google Cloud to manage and provide insights to their ever-growing datasets. The challenge of wrangling massive datasets is affecting almost every industry, including healthcare and life sciences, financial services, weather and climatology, and scientific research across a wide set of domains. In financial services, for example, historical market data can be stored in Google Cloud Storage. When developing new program trading algorithms, there is a need to back test (e.g., verify) against this data. HPDA applications can read input data from GCS, then use DAOS as a scratch space for intermediate data, running with high efficiency for hours or days, and then write their final output data back to GCS before

shutting down the scratch DAOS system. DAOS promises to speed processing of petabytes of data by removing I/O bottlenecks, and that means users can use fewer compute VMs and for a shorter period of time to achieve a more cost-effective and timely result.

## Synthetic Benchmark Results

> It is rare to see such good performance from a single storage system across all four dimensions.
>
> – Dean Hildebrand

Initial performance results on Google Cloud are very encouraging according to the key HPC metrics of MDop/s, IOP/s, bandwidth, and latency. The consensus by both teams, expressed by Dean Hildebrand, Technical Director in the Google Cloud Office of the CTO, is, "It is rare to see such good performance from a single storage system across all four dimensions."

The configuration for the IO500 benchmark runs are shown below (Figure 7).

| | |
|---|---|
| DAOS Clients | 32 x c2-standard-30 (Cascade Lake) with 120GB RAM, 32Gbps networking |
| DAOS Servers | 17 x x n2-custom-36 (Ice Lake) with 256GB RAM, 6TB Local SSD, 50Gbps networking |
| Storage Configuration | 102TB (raw), no erasure coding |

Figure 7. Benchmark hardware configuration.

The following table (Figure 8) reports the performance results.

| Bandwidth | |
|---|---|
| Read | 96 GiB/s (768 Tbps) |
| Write | 60 GiB/s (480 Tbps) |
| **Metadata OPs** | |
| Stat | 1M/sec |
| File create (empty) | 1.5M/sec |
| File create (3901 bytes) | 689K/sec |
| Small file read (3901 bytes) | 551K/sec |
| **IOPs and Latency** | |
| Random read (4KB) | 800K/sec (Latency 0.28ms) |
| Random write (4KB) | 825K/sec (Latency 0.36ms) |

Figure 8. Benchmark results.

For bandwidth, DAOS achieved extremely high efficiency, realizing over 94% of the published VM network and Local-SSD bandwidth. The differentiation of DAOS is through the very high IOP/s, MDop/s, and remarkably low latency even with the relatively small deployment tested. DAOS is able to read and write small I/O to large files at over 800K IOPS/sec, all with approximately 0.3ms latency at scale. Further, DAOS is able to process small files at 100s of thousands to millions

per second. The performance characteristics are critical as datasets in finance, life sciences, astronomy, and others can consist of billions of small files and large files with small pieces of information stored inside of them.

The benchmark in Figure 8 was optimized to reflect the use of DAOS as HPC scratch storage. For this reason, the DAOS configuration did not employ cross-VM redundancy. DAOS supports a robust fault model that ensures data and metadata integrity through the use of replication and erasure coding across VMs. For users that need increased resiliency, DAOS's software redundancy support (compared to traditional hardware support) provides the flexibility to dynamically grow and shrink storage capacity. Of course, additional resiliency will consume additional storage capacity and slightly lower write performance but read performance should be unaffected.

## TensorFlow-IO Performance for the CosmoFlow Application

During the SC'22 conference, Google and Intel teams demonstrated the performance of the CosmoFlow AI application leveraging the TensorFlow framework with DAOS. This reflects a common AI use case for the Google Cloud.

The CosmoFlow training application benchmark is part of the MLPerf HPC benchmark suite. It involves training a 3D convolutional neural network for N-body cosmology simulation data to predict physical parameters of the universe.[9] The benchmark built on top of the TensorFlow framework makes heavy use of convolution and pooling primitives. In total, the model contains approximately 8.9 million trainable parameters, and the model is trained on batches of large 3D images of approximately 3MB each. The resulting trained network is used to predict physical parameters of the universe.

The TensorFlow I/O (TF-IO) framework has been integrated with the DAOS libdfs I/O library (Figure 9), thus giving the application direct user space access to storage. As noted previously, this bypasses many POSIX and operating system kernel inefficiencies. No kernel modifications were required as the appropriate calls to libdfs were made in TF-IO. This reflects how direct user space integration can simplify HPC/AI/HPDA deployments in the cloud.
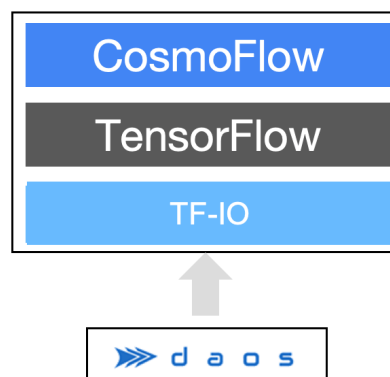


Figure 9. The TensorFlow-IO interface is built on top of the DAOS libdfs.

A preprocessed version of the dataset was used for the benchmark run.[11] The many individual files created for the benchmark are stored as individual files in the TFRecord format and provide an excellent stress test for the storage system.

The application is composed of a number of epochs. In each epoch, the records are loaded in parallel across the client ranks. In the standard CosmoFlow application these records would then be used for distributed training, but since this benchmark is focused on the I/O performance, only the data loading was performed (none of the distributed training phases were performed).

For the client side, the testing setup uses 29 c2-standard-60 nodes, where each node has 30 3.10GHz Intel® Xeon® processor cores, 240 GB RAM, and a gVNIC. For the server side, 16 n2-standard-80 nodes were used, each with 40 2.80GHz Intel Xeon processor cores, 640 GB RAM, a gVNIC, and 6 TB of Local SSD. The nodes all use TCP over ethernet and Tier 1 networking (100 Gbps). The benchmark ran for 305.7 seconds and performed 10 epochs, thus loading the entire dataset 10 times (5,242,890 files loaded, ~15.2 TB). DAOS achieved 49.7 GB/s, averaged loading 17,150 records per second, even as it repeatedly ingested a large number of small files.

## Building Working GCP Systems and Tuning

HPC/AI/HPDA workloads can leverage all the power of the Google Cloud for compute and data handling by incorporating DAOS with the Google HPC Toolkit.

The Google Cloud HPC Toolkit is open-source software offered by Google Cloud that makes it easy for customers to deploy HPC environments on Google Cloud. It allows customers to deploy turnkey HPC environments (compute, networking, storage, etc.) following Google Cloud best-practices, in a repeatable manner. The HPC Toolkit is designed to be highly customizable and extensible and intends to address the HPC deployment needs of a broad range of customers.

The Cloud HPC Toolkit reads a blueprint, which is a YAML file that defines a reusable configuration and describes the specific HPC environment that you want to deploy. The toolkit then generates deployments using Terraform, Packer and other widely used technologies.

The HPC Toolkit comes with several example configuration blueprints. These can be used as-is to get familiar with the operations of the HPC Toolkit, or they can be modified to build different configurations. Part of these examples are maintained by Intel and include the ability to deploy a cluster using Slurm, a DAOS cluster with client and server instances and a DAOS Server that is used by a Slurm cluster. While these examples are meant to demonstrate the ease of use of the Cloud HPC Toolkit, they can be customized and tuned to satisfy specific use cases.

## Generic, Extreme-Scale HPC/AI/HPDA Workload Support

Users can also choose from familiar HPC packages maintained outside of Google. For example, users can choose to run certified VM images on an Intel cluster. These workloads are optimized for HPC workloads on Intel hardware, satisfy the requirements for Intel Select Solution
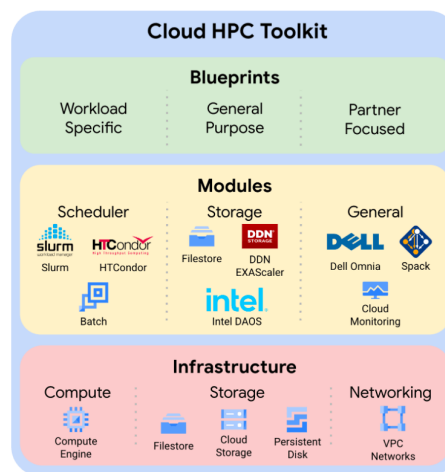


Figure 10. The Cloud HPC Toolkit provides many example blueprints that use various HPC Modules to create deployments on Google Cloud Infrastructure.

verification, and are verified compatible with applications listed in the Intel HPC application catalog.

For extreme scale-out, numerous concrete examples of HPC software can be found on the state-of-the-art Exascale Computing Project (ECP) applications and libraries. Examples include PETSc, SLATE, Ginkgo, heFFTe, SUNDIALS, hypre and more. The DOE-led Exascale Computing Initiative (ECI), a partnership between two DOE organizations, the Office of Science (SC) and the National Nuclear Security Administration (NNSA), was formed in 2016 to accelerate research, development, acquisition, and deployment projects to deliver exascale computing capability to the DOE laboratories by the early to mid-2020s. This includes leveraging GPU-acceleration and next-generation storage I/O such as DAOS, which will be the primary storage software for the Aurora exascale supercomputer to be sited at Argonne National Laboratory.

**Key points of the ECP software include:**

- The ability to download and run pre-built packages from the Extreme-scale Scientific Software Stack (E4S). Alternatively, build with Spack to run C++, SYCL, Fortran, OpenMP, MPI and Python source code. Intel notes that with integrated Intel® oneAPI tools support, these languages deliver productive cross-architecture performance across CPUs and GPUs using a single codebase.[12]

- Utilize the latest HPC applications and libraries including the Exascale Computing Project (ECP) projects. Given the Aurora supercomputer reliance on DAOS for IO, its compatibility is baked in by design and regularly tested with continuous integration.

- Visualize even the largest HPC data sets using situ visualization capabilities, even on your laptop and with Jupyter notebooks. Visualizations can even include ray-traced photorealistic images using software define visualization (SDVis) with the Intel® OSPRay library.[13]

From Google's pre-tuned HPC images to third-party images, the Google HPC Toolkit gives HPC users an automated experience and the ability to transition between on-premises and cloud supercomputing, as desired, and at a scale of the users choosing.

## Google C3 – 4th Gen Intel® Xeon® Processors Plus Advanced Benefits HPC Storage

The Google C3 machine instances, which are just starting to come online, provide access to the latest 4th Gen Intel Xeon processors. These Google Cloud C3 instances also provide hardware accelerated communications, providing low latency and accelerated and secure networking that benefits storage I/O. DAOS can build upon this hardware capability to access storage in the Google Cloud. This performance addition is available on the new C3 virtual machines (VMs), which utilizes a system-on-chip (SoC) design that couples a 4th Gen Intel Xeon CPU with a custom-built Intel infrastructure processing unit (IPU).

C3 VMs are designed to help cloud and network providers free up CPU resources by offloading functions like storage and network virtualization, to increase overall workload performance and also yields various security benefits. Nick McKeown, senior vice president and general manager of Network and Edge at Intel notes "A first of its kind in any public cloud, C3 VMs will run workloads on 4th Gen Intel Xeon processors while they free up programmable packet processing to the IPUs securely at line rates of 200Gb/s. This Intel and Google collaboration enables customers through infrastructure that is more secure, flexible, and performant."[14]

## Getting Started

Google Cloud is a great place to run with the latest hardware in a balanced and performant HPC-in-the-cloud environment. See this DAOS website for HPC Toolkit and Terraform deployment tooling and instructions to ensure your applications realize extreme storage performance, commensurate with Google's high compute and networking capabilities.

intel.

1 https://www.techradar.com/news/new-google-cloud-vms-will-be-the-first-to-run-on-intel-xeon-sapphire-rapids
2 https://www.hpcwire.com/2022/10/17/daos-performance-expands-beyond-intel-optane-and-into-the-google-cloud/
3 https://newsroom.intel.com/articles/intel-optane-persistent-memory-daos-solution-sets-world-record/#gs.dxi98n
4 https://www.alcf.anl.gov/support-center/aurora/aurora-introduction
5 https://community.intel.com/t5/Blogs/Products-and-Solutions/HPC/DAOS-Momentum-Demonstrated-with-New-IO500-Rankings-and-Community/post/1389619
6 Running in user space eliminates the overhead of calling the kernel to perform a file operation, which requires a syscall and a context switch. Direct user space access also eliminates the overhead of servicing the operating system virtual file system page and block caches. Also see https://link.springer.com/chapter/10.1007/978-3-030-48842-0_3.
7 DAOS uses a key-value architecture that eliminates many POSIX performance and scalability issues. Also see https://medium.com/@rmfarber/hurray-for-daos-exit-the-bottlenecks-imposed-by-posix-io-479feba87ec1
8 https://www.intel.com/content/www/us/en/high-performance-computing/daos.html
9 https://proxyapps.exascaleproject.org/app/mlperf-cosmoflow/
10 https://arxiv.org/abs/1808.04728
11 https://github.com/sparticlesteve/cosmoflow-benchmark
12 https://www.hpcwire.com/off-the-wire/intel-and-google-cloud-announce-cloud-hpc-toolkit/
13 https://www.exascaleproject.org/highlight/visualization-and-analysis-with-cinema-in-the-exascale-era/
14 https://www.hpcwire.com/off-the-wire/intel-and-google-cloud-announce-cloud-hpc-toolkit/