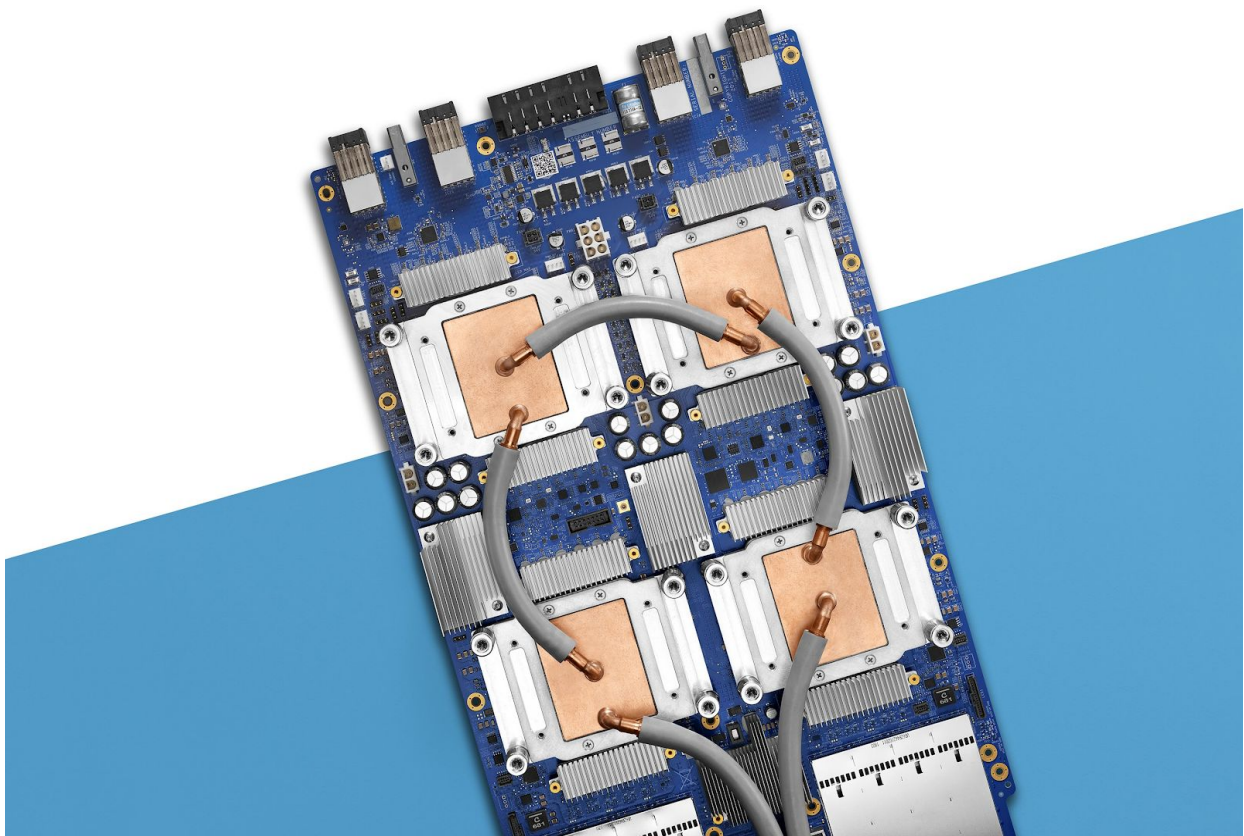


High Performance Computing (HPC) Solutions in the Cloud

Cloud training, resources, and reference architectures for researchers and data scientists.



Google Cloud TPU v3

Contents

Introduction	2
Learning	2
High Performance Computing Solutions	4
Singularity Containers with Cloud Build	5
Host Jupyter Notebooks in a Slurm Cluster	7
GKE Dynamic Storage Provisioning from Cloud Filestore	9
TensorFlow Inference Workloads at Scale with TensorRT 5	11
Appendix	13

Introduction

The National Institutes of Health (NIH) established The Science and Technology Research Infrastructure for Discovery, Experimentation, and Sustainability (STRIDES) initiative to provide biomedical researchers with access to advanced, cost-effective, cloud-based computational infrastructure, tools, and services. Through STRIDES, researchers can take advantage of emerging data management methodologies, technological expertise, computational platforms, and tools to support cutting-edge experimentation and innovation. NIH has partnered with Google Cloud to support the STRIDES initiative through cloud services. In support of STRIDES, we've developed sets of playbooks to help enable researchers to build healthcare and life sciences solutions on Google Cloud Platform (GCP).

The goal of this playbook is to assist researchers in designing and implementing High Performance Computing workflows on GCP. More specifically, this playbook will focus on architecture and solutions examples for **building and storing singularity containers, hosting Jupyter notebooks in slurm clusters, setting up NFS file storage for Google Kubernetes Engine (GKE), and running TensorFlow inference workloads at scale.** Additionally, this playbook will outline training and digital resources to help upskill and enable researchers to build on Google Cloud, while highlighting the appropriate products and services to use when architecting on GCP.

Learning

Generally, cloud adopters fall under one of three categories:

Cloud Novice	Cloud Ready	Cloud Native
Little to no experience with the cloud	Familiar with the cloud, some experience	Lots of cloud experience, expert-level knowledge

Understanding this broad spectrum of experience levels, we've highlighted key training resources to help upskill researchers on Google Cloud. Additionally, Google offers [on-site, instructor-led training](#) to enable large groups of participants across your organization.

Google Cloud

Cloud Novice

[Video: Welcome to GCP](#)

[Documentation: GCP Conceptual Overview](#)

[Documentation: All GCP Products & Services](#)

[Video: Intro to GCP for Students](#)

[Documentation: About GCP Services](#)

[Virtual Course: GCP Fundamentals - Core Infrastructure](#)

[Video: GCP 101](#)

[Documentation: GCP Development & Admin Tools](#)

[Virtual Lab: GCP Essentials](#)

[Video: GCP Essentials](#)

Cloud Ready

[Documentation: All GCP Products & Services](#)

[Video: Deploy an HPC Cluster](#)

[Infrastructure - Scaling and Automation](#)

[Documentation: GCP for Data Center Professionals](#)

[Video: Intro to HPC on GCP](#)

[Virtual Course: Elastic Cloud Infrastructure - Containers and Services](#)

[Documentation: GCP for AWS Professionals](#)

[Video: GPUs on GCP](#)

[Virtual Lab: Cloud Architecture](#)

[Documentation: GCP for Azure Professionals](#)

[Virtual Course: Essential Cloud Infrastructure - Foundation](#)

[Virtual Lab: Cloud Engineering](#)

[Documentation: GCP for OpenStack Users](#)

[Virtual Course: Essential Cloud Infrastructure - Core Services](#)
[Virtual Course: Elastic Cloud](#)

[Virtual Lab: Cloud Development](#)

[Virtual Course: Architecting with Google Cloud Platform](#)

Cloud Native

[Documentation: All GCP Products & Services](#)

[Video: Storage for HPC](#)

[Machine Learning Fundamentals](#)

[Documentation: GCP Solutions](#)

[Video: GPU Infrastructure for ML and HPC Workloads](#)

[Virtual Course: Leveraging Unstructured Data with Cloud Dataproc on Google Cloud Platform](#)

[Documentation: Cloud Tensor Processing Units \(TPUs\)](#)

[Tutorial: Running R at Scale](#)

[Virtual Course: Serverless Data Analysis with Google BigQuery and Cloud Dataflow](#)

[Code Sample: HT Condor Deployment Manager Scripts](#)

[Virtual Course: Machine Learning with TensorFlow](#)

[Codelab: Genomics Batch and High-Throughput Computing](#)

[Virtual Lab: Intermediate ML: Tensorflow on GCP](#)

[Virtual Course: Serverless Machine Learning with Tensorflow on Google Cloud Platform](#)

[Code Sample: Lustre Parallel File System Deployment Scripts](#)

[Codelab: Image Processing & Astronomy](#)

[Virtual Course: Big Data and](#)

[Virtual Course: Building Resilient Streaming Systems on GCP](#)

[Virtual Lab: Data Engineering](#)

[Virtual Lab: Data Science on GCP](#)

[Virtual Lab: Data Science on GCP - Machine Learning](#)

[Tutorial: Genomics GATK](#)

[Pipeline Virtual Lab: Google Cloud Solutions II - Data and ML](#)

[Tutorial: TensorFlow Inference at Scale](#)

High Performance Computing Solutions

Google offers powerful, flexible infrastructure capable of supporting scalable workloads. This allows researchers to work with data at a larger scale and faster than before, decreasing time from hypothesis to insight. GCP can enable and accelerate time to results for tasks from analyzing genomic sequence data to simulating complex climate models.

Compute, storage, and networking all operate at global scale at Google. Whether you need petabytes of storage or the latest processors, Google's infrastructure delivers. Build customized clusters with dozens of CPUs, powerful [NVIDIA GPUs](#) and [Cloud Tensor Processing Units](#) (TPUs), and high-throughput and low-latency [object](#) and [file storage](#).

Cloud HPC solutions mean that every team has access to an effectively infinitely-scalable setup, tailored to their specific needs. This means less time waiting for workloads in queues, and more time for the work that really matters. For batch workloads, GCP offers the cloud-native [Batch](#), providing dynamically-allocated compute clusters on [Google Kubernetes Engine](#) (GKE).

Google Cloud compute resources are billed at a per-second rate, meaning that you only pay for what you use. There are also a number of discounts available, such as automatically-applied Sustained Use Discounts, and up to 80% cost-savings with [Preemptible VMs](#).

Tools such as [BigQuery](#), [AI Platform](#) and [Tensorflow](#) enable researchers to apply analytics and artificial intelligence to data. Together with [AutoML](#), users can establish mechanisms to auto-detect patterns, predict outcomes, and quickly analyze large amounts of data.

Additionally, both G Suite and GCP are [FedRAMP compliant](#) and [HITRUST CSF certified](#); prescriptive security controls are implemented throughout Google Cloud for processing, storing, and transmitting sensitive data. Both [G Suite](#) and [GCP](#) also support HIPAA compliance across dozens of products in all regions, zones, network paths, and points of presence for Google Cloud.

Google Cloud

[Learn more](#) about Google Cloud's high performance computing and innovative capabilities.

Reference Architectures

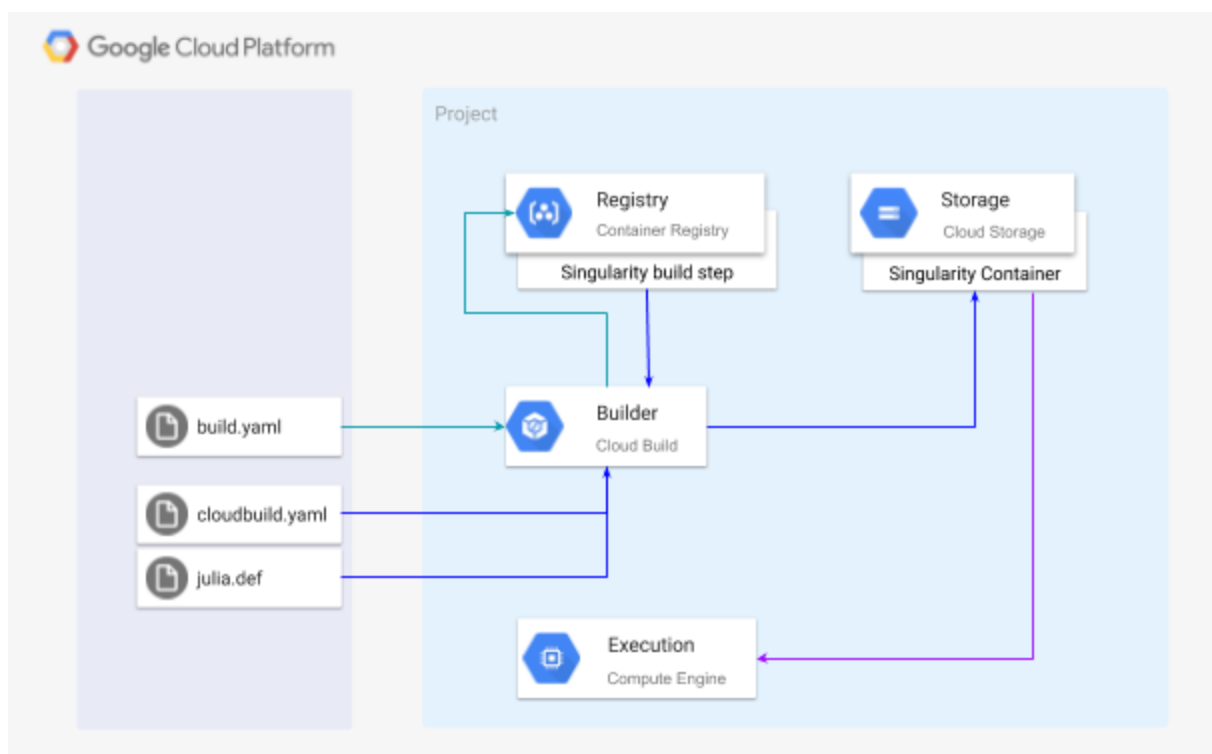
[\(link\)](#)

Google Cloud provides access to vastly scalable infrastructure backed by cutting-edge hardware to support your high performance workflows. Included here are a number of common HPC architectures and guidance on how to get them set up.

Singularity Containers with Cloud Build

[\(link\)](#)

Build, store, and deploy HPC-optimized [Singularity](#) containers on GCP with [Cloud Build](#).



Solution Summary:

Use the cloud to configure and manage Singularity containers. Unlike Docker, the Singularity binary is specifically designed for HPC applications. Cloud Build provides a standard set of build steps which you can extend by creating your own specialized images. In this solution, we define a custom build step to build Singularity container images. We then use this to construct a Singularity container artifact and store it in Cloud Storage, before deploying to Compute Engine.

Google Cloud

Suggested GCP Products and Services:

- [Cloud VPC Network](#) - Virtual Private Cloud provides global, scalable, and flexible networking for your GCP resources and services without compromising integrity on public-facing connections.
- [Cloud Build](#) - Executes your container builds on GCP infrastructure. Builds are run as a sequence of build steps, each within a Docker container. You can use provided build steps, or define your own custom ones.
- [Container Registry](#) - Private container image registry that runs on Google Cloud. Supports Docker Image Manifest V2 and OCI image formats.
- [Cloud Storage](#) - Unified object storage for storing raw data files and processed data. Supports regional, multi-regional, archive, and infrequently accessed data.
- [Google Compute Engine \(GCE\)](#) - VM instances for processing data.

High-level Setup Steps:

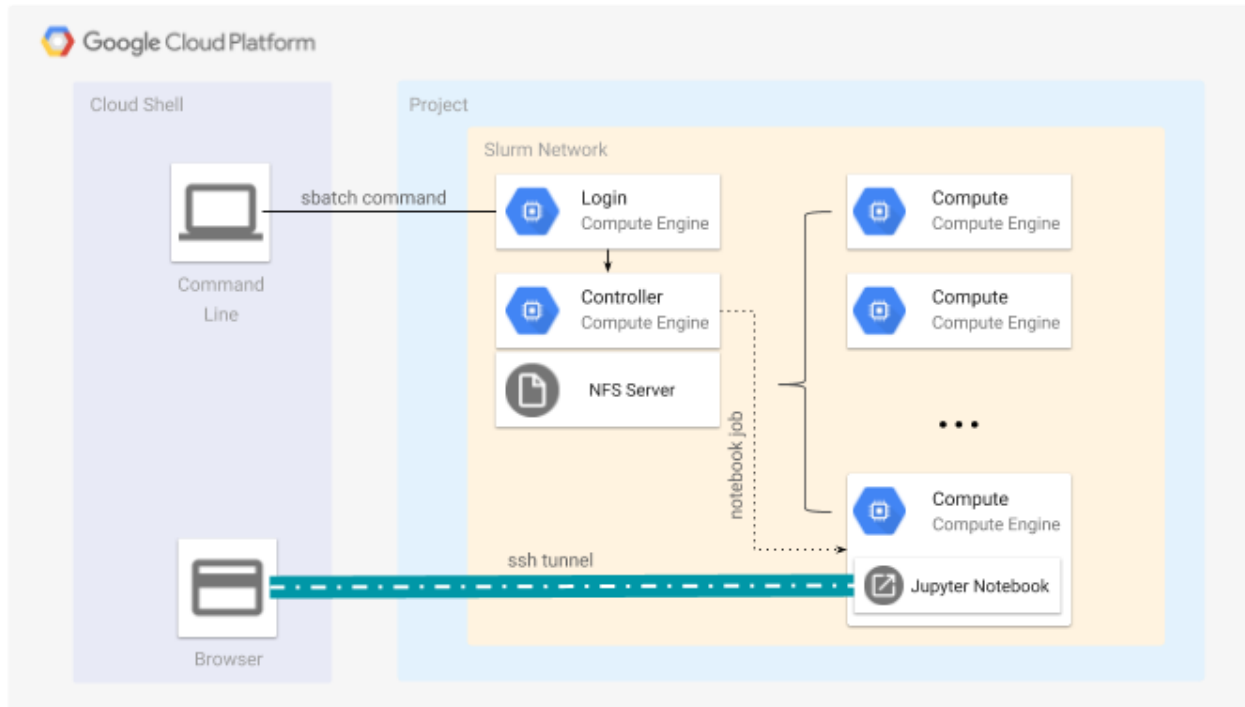
1. Create a [Google Cloud Project](#) and use [Cloud IAM](#) to manage who has access to the resources within the GCP project.
2. Enable the Compute Engine, Cloud Build, and Container Registry APIs.
3. Define, and then build, a custom build step in `build.yaml` to create the Singularity builder.
4. Create a [Cloud Storage bucket](#) that will store the built container artifact. Note that you can use [IAM](#) to configure [access policies](#) to govern which workloads can access the containers you build.
5. Create a Singularity definition file (`julia.def`) to establish a blueprint for how the container will be constructed.
6. In `cloudbuild.yaml`, specify the steps required to build the container and store the result in a Cloud Storage bucket.
7. Use the custom Singularity build step created earlier together with `cloudbuild.yaml` and `julia.def` to build and store the container.
8. To deploy the container, first create a [Compute Engine virtual machine](#) with the singularity binary installed. Then, [download the container from Cloud Storage](#) and run it using `singularity`.

For code samples and a more detailed walkthrough, refer [here](#).

Host Jupyter Notebooks in a Slurm Cluster

(link)

Run a [Jupyter Notebook](#) as a job managed by [Slurm Workload Manager](#).



Solution Summary:

Slurm is a popular resource manager used in many high performance computing centers. Jupyter Notebooks are a favorite tool of machine learning and data science specialists. There are situations that call for specialized hardware, such as GPUs, or more memory or CPU cores than are available locally. In these situations, Slurm can allocate compute instances that have the required hardware resources to run the user's notebook for a bounded time period.

Suggested GCP Products and Services:

- [Google Compute Engine \(GCE\)](#) - VM instances for processing data.
- [Deployment Manager](#) - Infrastructure deployment service that automates the creation and management of Google Cloud resources.
- [Cloud VPC Network](#) - Virtual Private Cloud provides global, scalable, and flexible networking for your GCP resources and services without compromising integrity on public-facing connections.

High-level Setup Steps:

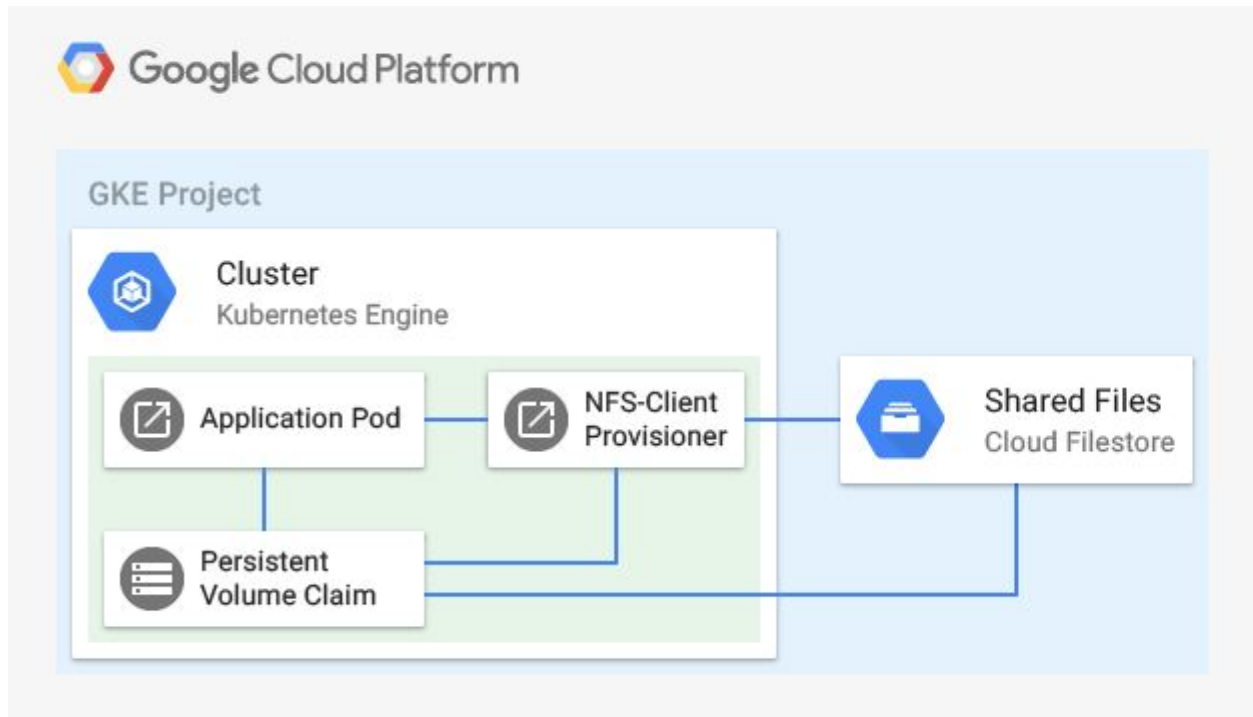
1. Create a [Google Cloud Project](#) and use [Cloud IAM](#) to manage who has access to the resources within the GCP project.
2. Enable the Compute Engine API.
3. Clone the [Slurm on Google Cloud Platform](#) GitHub repository, and customize the configuration for your environment.
4. Use [Deployment Manager](#) to deploy the Slurm Resource Manager to a Compute Engine cluster.
5. Setup an [Anaconda](#) environment module by installing the package on the controller, where each Slurm compute node can access it.
6. Run a Jupyter Notebook as a batch job on the Slurm cluster. This is done by submitting a notebook .batch file to Slurm via the sbatch command.
7. Establish an SSH tunnel to your Notebook via the Slurm login node. Use your browser to connect to your Notebook's interface through this tunnel. You now have a Jupyter Notebook up and running.

For code samples and a more detailed walkthrough, refer [here](#).

GKE Dynamic Storage Provisioning from Cloud Filestore

([link](#))

Dynamically provision storage volumes in [Google Kubernetes Engine](#) (GKE) from [Cloud Filestore](#) using an [NFS-Client Provisioner](#).



Solution Summary:

Dynamic provisioning allows storage volumes to be created on demand in an NFS volume managed by a Cloud Filestore instance. The NFS-Client Provisioner should be deployed to fulfil persistent volume storage claims from the application pods using Cloud Filestore's managed storage. Typical use cases include running databases, e.g., PostgreSQL, or a content management system like WordPress in a Kubernetes cluster.

Suggested GCP Products and Services:

- [Google Kubernetes Engine](#) - GKE provides a managed environment for deploying, managing, and scaling your containerized applications.
- [Cloud Filestore](#) - Managed file storage service for applications that require a filesystem interface and a shared file system for data, such as Network Attached Storage (NAS).

Google Cloud

- [Cloud VPC Network](#) - Virtual Private Cloud provides global, scalable, and flexible networking for your GCP resources and services without compromising integrity on public-facing connections.

High-level Setup Steps:

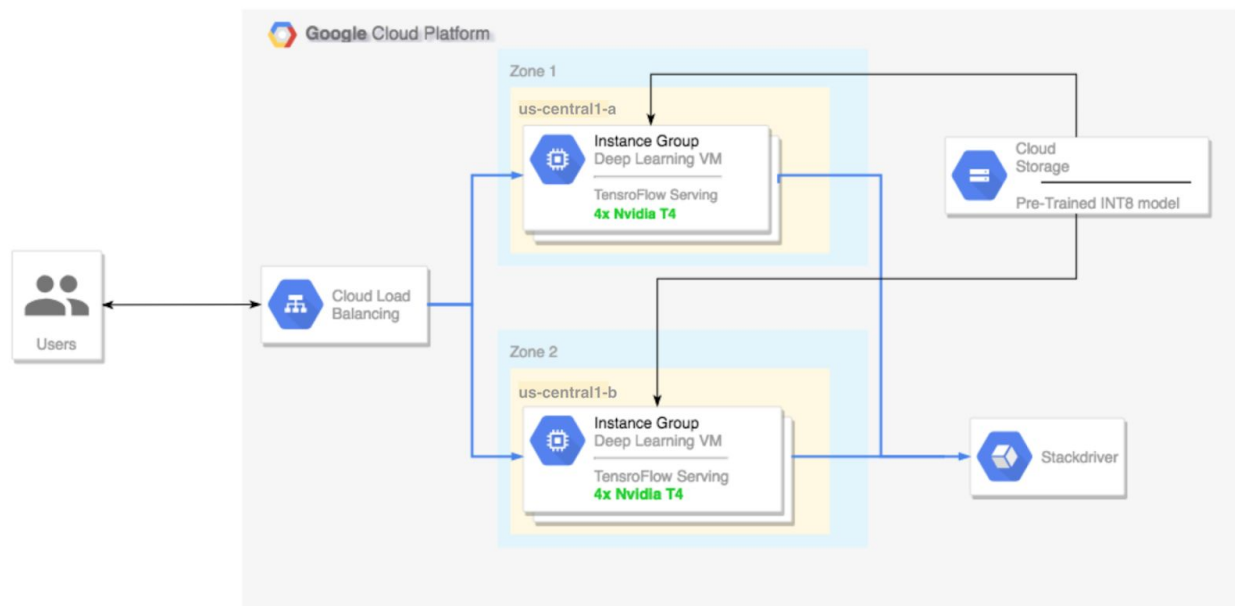
1. Create a [Google Cloud Project](#) and use [Cloud IAM](#) to manage who has access to the resources within the GCP project.
2. In the GCP project, [create a VPC network](#) to logically isolate your project resources.
3. Ensure the Cloud Filestore and GKE APIs are enabled.
4. Set up a Cloud Filestore storage volume.
5. Deploy a [GKE cluster](#). This will host your application pods, and also the NFS-Client Provisioner.
6. Install the NFS-Client Provisioner using [Helm](#) and deploy it to the cluster.
7. You are up and running! Make a [Persistent Volume Claim](#) using any application that uses storage classes to do dynamic provisioning.
8. Optional: Verify Cloud Filestore volume directory creation by mounting the volume on a Compute Engine instance and inspecting the directory structure to see that the database files are present.

For code samples and a more detailed walkthrough, refer [here](#).

TensorFlow Inference Workloads at Scale with TensorRT 5

([link](#))

Run deep learning inference on large-scale workloads with [NVIDIA TensorRT 5](#) and [NVIDIA T4](#) GPUs.



Solution Summary:

Set up a multi-zone cluster for running an inference workload on an autoscaling group that scales to meet changing GPU utilization demands. The Compute Engine VMs will be configured with Cloud Deep Learning VM images, and will serve the model using TensorFlow Serving.

Suggested GCP Products and Services:

- [Google Compute Engine \(GCE\)](#) - VM instances for processing data.
- [Cloud Load Balancing](#) - Distribute load-balanced compute resources in single or multiple regions, meet high availability requirements, put resources behind a single anycast IP and scale resources up or down with intelligent Autoscaling.
- [Cloud Storage](#) - Unified object storage for storing raw data files and processed data. Supports regional, multi-regional, archive, and infrequently accessed data.
- [Stackdriver Monitoring](#) - Collects metrics, events, and metadata from GCP, AWS, hosted uptime probes, and application instrumentation, as well as over 150 common application components, on-premise systems, and hybrid cloud systems.

Google Cloud

- [Cloud VPC Network](#) - Virtual Private Cloud provides global, scalable, and flexible networking for your GCP resources and services without compromising integrity on public-facing connections.

High-level Setup Steps:

1. Create a Compute Engine VM configured with a [TensorFlow Deep Learning image](#) and in any zone that supports T4 GPUs.
2. Download and extract your TensorFlow [SavedModel](#), before converting it to a TensorRT graph with [TFTools](#).
3. Archive and upload the model to Cloud Storage, so that it can be served from the Compute Engine cluster later on.
4. Set up a multi-zone Compute Engine cluster:
 - a. Prepare a [startup script](#) to install required services and libraries, such as NVIDIA drivers, [GPU Utilization Agent](#), and [TensorFlow Serving](#).
 - b. Define an [Instance Template](#).
 - c. Create a [Managed Instance Group](#).
 - d. Configure [Autoscaling](#) using the GPU Utilization metric in Stackdriver.
5. Configure a Load Balancer in front of the Managed Instance Group. Ensure that [firewall rules](#) are set up to allow [health checks](#) by the Load Balancer, as well as HTTP traffic.
6. You can now POST model inputs to the TensorFlow Serving server using the Load Balancer's IP address to run GPU-backed inference and receive a prediction.

For code samples and a more detailed walkthrough, refer [here](#).

Appendix

High Performance Computing Solutions Reference Architectures [\(link\)](#)

Technical Resources for HPC [\(link\)](#)

Reference Architecture for Creating a HIPAA-Aligned Project in GCP [\(link\)](#)

High Performance Computing - Google Cloud Blog [\(link\)](#)

Using Clusters for Large-Scale Technical Computing in the Cloud [\(link\)](#)

GCP Technical Tutorials [\(link\)](#)