



Deploying Applications to Google Cloud

Learning objectives

- Choose the appropriate Google Cloud deployment service for your applications.
- Configure scalable, resilient infrastructure using Instance Groups.
- Orchestrate microservice deployments using Kubernetes, GKE and Cloud Run.
- Leverage App Engine for a completely automated platform as a service (PaaS).
- Create serverless applications using Cloud Functions.



Google Cloud offers several kinds of compute resources for deployment. Each offers different levels of control and features. In this module, we will discuss, compare, and contrast the following five services:

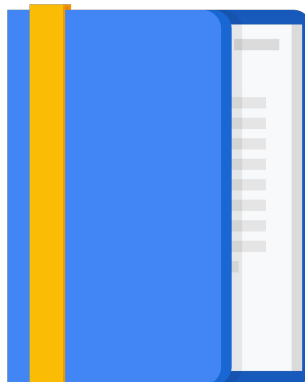
- Compute Engine
- App Engine
- Google Kubernetes Engine
- Cloud Functions
- Cloud Run

Agenda

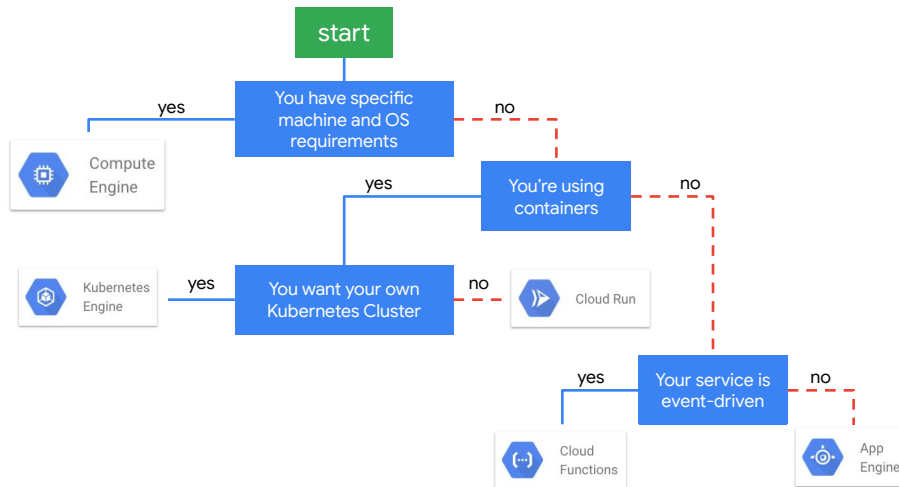
Google Cloud Infrastructure as a Service

Google Cloud Deployment Platforms

Lab



Choosing a Google Cloud deployment platform



<script>

Let me give you a high-level overview of how you could decide on the most suitable platform for your application.

First, ask yourself whether you have specific machine and OS requirements. If you do, then Compute Engine is the platform of choice.

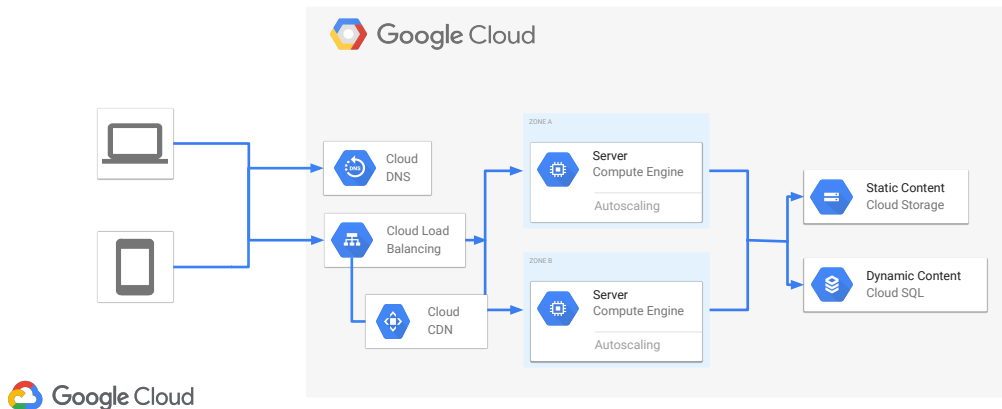
If you have no specific machine or operating system requirements, then the next question to ask is whether you are using containers. If you are, then you should consider Google Kubernetes Engine or Cloud Run, depending on whether you want to configure your own Kubernetes cluster.

If you are not using containers, then you want to consider Cloud Functions if your service is event-driven and App Engine if it's not.

We'll talk through each of these services in this module and you will get to explore them in a lab.

</script>

Use Compute Engine when you need complete control over operating systems, for apps that are not containerized or self-hosted databases



<script>

Compute Engine is a great solution when you need complete control over your operating systems, or if you have an application that is not containerized, an application built on a microservice architecture, or an application that is a database.

Instance groups and autoscaling as shown on this slide allow you to meet variations in demand on your application. Let's take a closer look at instance groups.

</script>

Managed instance groups create VMs based on instance templates

- Instance templates define the VMs: image, machine type, etc.
 - Test to find the smallest machine type that will run your program.
 - Use a Startup Script to install your program from a Git repo.
- Instance group manager creates the machines.
 - Set up auto scaling to optimize cost and meet varying user workloads.
 - Add a health check to enable auto healing.
 - Use multiple zones for high availability.



<script>

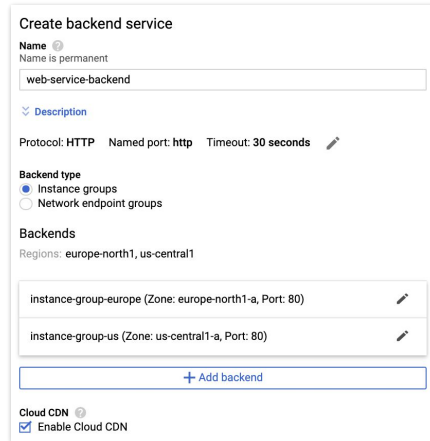
Managed instance groups create VMs based on instance templates. Instance templates are just a resource used to define VMs and managed instance groups. The templates define the boot disk image or container image to be used, the machine type, labels, and other instance properties like a startup script to install software from a Git repository.

The virtual machines in a managed instance group are created by an instance group manager. Using a managed instance group offers many advantages, such as autohealing to re-create instances that don't respond and creating instances in multiple zones for high availability.

</script>

Use one or more instance groups as the backend for load balancers

- Use a global load balancer if you have instance groups in multiple regions.
- Enable the CDN to cache static content.
- For external services, set up SSL.
- For internal services, don't provide a public IP address.



The screenshot shows the 'Create backend service' configuration page. The 'Name' field is set to 'web-service-backend'. The 'Description' section shows 'Protocol: HTTP', 'Named port: http', and 'Timeout: 30 seconds'. Under 'Backend type', 'Instance groups' is selected. The 'Backends' section lists two instance groups: 'instance-group-europe (Zone: europe-north1-a, Port: 80)' and 'instance-group-us (Zone: us-central1-a, Port: 80)'. At the bottom, the 'Cloud CDN' section has the 'Enable Cloud CDN' checkbox checked.

Create backend service

Name ⓘ
Name is permanent
web-service-backend

Description ⓘ
Protocol: HTTP Named port: http Timeout: 30 seconds ✎

Backend type
☒ Instance groups
☐ Network endpoint groups

Backends
Regions: europe-north1, us-central1

instance-group-europe (Zone: europe-north1-a, Port: 80)	✎
instance-group-us (Zone: us-central1-a, Port: 80)	✎

+ Add backend

Cloud CDN ⓘ
☒ Enable Cloud CDN



<script>

I recommend using one or more instance groups as the backend for load balancers. If you need instance groups in multiple regions, use a global load balancer, and if you have static content, simply enable Cloud CDN as shown on the right.

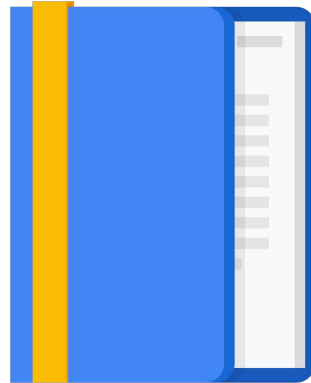
</script>

Agenda

Google Cloud Infrastructure as a Service

Google Cloud Deployment Platforms

Lab



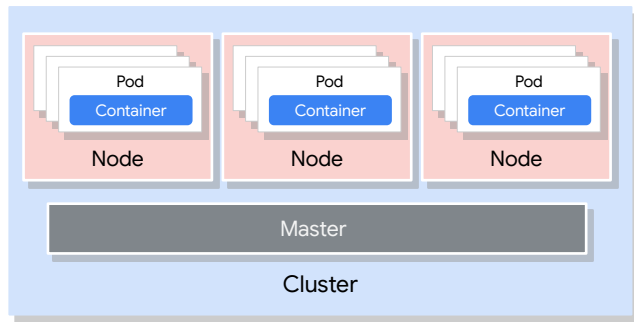
<script>

Let's go through the other deployment platforms: GKE, Cloud Run, App Engine, and Cloud Functions.

</script>

Google Kubernetes Engine (GKE) automates the creation and management of compute infrastructure

- Kubernetes clusters have a collection of nodes.
- In GKE, nodes are Compute Engine VMs.
- Services are deployed into pods.
- Optimize resource utilization by deploying multiple services to the same cluster.
- You pay for the VMs.



<script>

Google Kubernetes Engine, or GKE, provides a managed environment for deploying, managing, and scaling containerized applications using Google infrastructure. The GKE environment consists of multiple Compute Engine virtual machines grouped together to form a cluster. GKE clusters are powered by the Kubernetes open source cluster management system. Kubernetes provides the mechanisms with which to interact with the cluster. Kubernetes commands and resources are used to deploy and manage applications, perform administration tasks and set policies, and monitor the health of deployed workloads.

This diagram on the right shows the layout of a Kubernetes cluster. A cluster consists of at least one cluster master and multiple worker machines that are called nodes. These master and node machines run the Kubernetes cluster orchestration system. Pods are the smallest, most basic deployable objects in Kubernetes. A pod represents a single instance of a running process in a cluster. Pods contain one or more containers, such as Docker containers, that run the services being deployed. You can optimize resource use by deploying multiple services to the same cluster.

</script>

Cloud Run allows you to deploy containers to Google managed Kubernetes clusters

- Cloud Run allows you to use Kubernetes without the cluster management or configuration code.
- Apps must be stateless.
- Need to deploy apps using Docker images in Container Registry.
- Can also use Cloud Run to automate deployment to your own GKE cluster.

Container

Container image URL *

[SELECT](#)

E.g. gcr.io/cloudrun/hello

Must be stateless and listen for HTTP requests on \$PORT. [How to build a container?](#)

Deployment platform ⓘ

☒ Cloud Run (fully managed)

Location *
 ▼

Region for this Service can't be changed later. [How to pick a region?](#)

☐ Cloud Run for Anthos



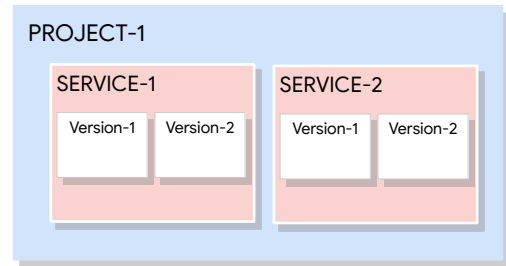
<script>

Cloud Run, on the other hand, allows you to deploy containers to a Google-managed Kubernetes cluster. A big advantage is that you don't need to manage or configure the cluster. The services that you deploy must be stateless, and the images you use must be in Container Registry. Cloud Run can be used to automate deployment to Anthos GKE clusters. You should do this if you need more control over your services, because it will allow you to access your VPC network, tune the size of compute instances, and run your services in all GKE regions. The screenshot on the right shows a Cloud Run configuration where the container image URL is specified along with the deployment platform, which can be fully managed Cloud Run or Cloud Run for Anthos.

</script>

App Engine was designed for microservices

- Each Google Cloud project can contain 1 App Engine application.
- An application has 1 or more services.
- Each service has 1 or more versions.
- Versions have 1 or more instances.
- Automatic traffic splitting for switching versions.



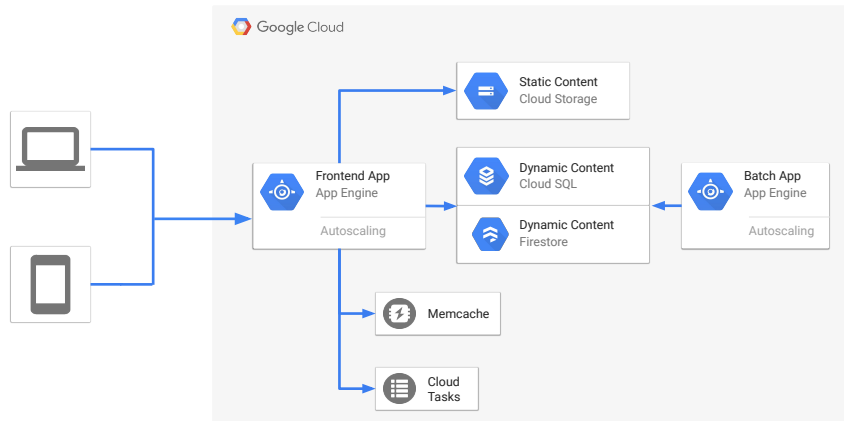
<script>

App Engine is a fully managed, serverless application platform supporting the building and deploying of applications. Applications can be scaled seamlessly from zero upward without having to worry about managing the underlying infrastructure. App Engine was designed for microservices. For configuration, each Google Cloud project can contain one App Engine application, and an application has one or more services. Each service can have one or more versions, and each version has one or more instances. App Engine supports traffic splitting so it makes switching between versions and strategies such as canary testing or A/B testing simple. The diagram on the right shows the high-level organization of a Google Cloud project with two services, and each service has two versions. These services are independently deployable and versioned.

Let me show you a typical App Engine microservice architecture.

</script>

Typical App Engine microservice architecture



<script>

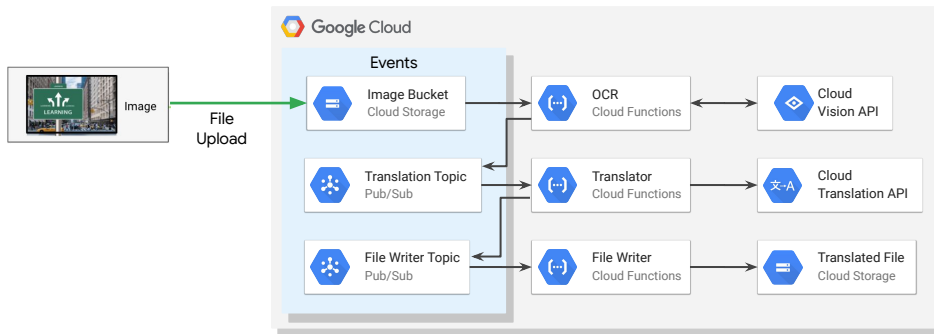
This could be an example of a retailer that sells online. Here App Engine serves as the frontend for both web and mobile clients. The backend of this application is a variety of Google Cloud storage solutions with static content such as images stored in Cloud Storage, Cloud SQL used for structured relational data such as customer data and sales data, and Firestore used for NoSQL storage such as product data. Firestore has the benefit of being able to synchronize with client applications.

Memcache is used to reduce the load on the datastores by caching queries, and Cloud Tasks are used to perform work asynchronously outside a user request (or service-service request). There's also a batch application that generates data reports for management.

</script>

Cloud Functions is great way to create loosely coupled, event-driven microservices

- Can be triggered by changes in a storage bucket, Pub/Sub messages, or web requests
- Completely managed, scalable, and inexpensive



<script>

Cloud Functions are a great way to deploy loosely coupled, event-driven microservices. They have been designed for processing events that occur in Google Cloud. The functions can be triggered by changes in a Cloud Storage bucket, a Pub/Sub message, or HTTP requests. The platform is completely managed, scalable, and inexpensive. You do not pay if there are no requests, and processing is paid for by execution time in 100ms increments.

This graphic illustrates an image translation service implemented with Cloud Functions. When an image is uploaded to a Cloud Storage bucket, it triggers an OCR Cloud Function that identifies the text in the image using Google's Cloud Vision API. Once the text has been identified, this service then publishes a message to a Pub/Sub topic for translation, which triggers another Cloud Function that will translate the identified text in the image using the Cloud Translation API. After that, the translator Cloud Function will publish a message to a file write topic in Pub/Sub, which triggers a Cloud Function that will write the translated image to a file.

This sequence illustrates a typical use case of Cloud Functions for event-based processing.

</script>

Lab

Deploying Apps to Google Cloud



App Engine



Google
Kubernetes
Engine



Cloud Run

Objectives

- Deploy to App Engine
- Deploy to Google Kubernetes Engine
- Deploy to Cloud Run

Quiz

You need to deploy an existing application that was written in .NET version 4. The application requires Windows servers, and you don't want to change it. Which should you use?

- A. Compute Engine
- B. GKE
- C. App Engine
- D. Cloud Functions



You need to deploy an existing application that was written in .NET version 4. The application requires Windows servers, and you don't want to change it. Which should you use?

- A. Compute Engine
- B. GKE
- C. App Engine
- D. Cloud Functions

Quiz

You need to deploy an existing application that was written in .NET version 4. The application requires Windows servers, and you don't want to change it. Which should you use?

A. Compute Engine

B. GKE

C. App Engine

D. Cloud Functions



- A. This is the correct answer. The approach is a lift and shift which is best supported by Compute Engine, because Compute Engine offers full control over virtual machines including operating system. No repackaging would be required.
- B. This answer is not correct. GKE would require repackaging into Docker containers.
- C. This answer is not correct. App Engine standard environment does not support .NET, and using App Engine flexible environment would require repackaging the application.
- D. This answer is not correct. Cloud Functions is the wrong model; it is for deploying single purpose functions.

Quiz

You have containerized multiple applications using Docker and have deployed them using Compute Engine VMs. You want to save costs and simplify container management. What might you do?

- A. Write Terraform scripts for all deployment.
- B. Rewrite the applications to run in App Engine.
- C. Rewrite the applications to run in Cloud Functions.
- D. Migrate the containers to GKE.



You have containerized multiple applications using Docker and have deployed them using Compute Engine VMs. You want to save costs and simplify container management. What might you do?

- A. Write Terraform scripts for all deployment.
- B. Rewrite the applications to run in App Engine.
- C. Rewrite the applications to run in Cloud Functions.
- D. Migrate the containers to GKE.

Quiz

You have containerized multiple applications using Docker and have deployed them using Compute Engine VMs. You want to save costs and simplify container management. What might you do?

- A. Write Terraform scripts for all deployment.
- B. Rewrite the applications to run in App Engine.
- C. Rewrite the applications to run in Cloud Functions.
- D. Migrate the containers to GKE.



- A. This answer is not correct. While this could be achieved with Terraform, the requirement to save costs and simplify container management would not be met.
- B. If the applications are containerized, then rewriting to run in App Engine is not cost-effective.
- C. This answer is not correct. Cloud Functions are for deploying single purpose functions, not applications.
- D. This is the correct answer. The applications are containerized, and GKE will help with the resource efficiency and hence costs, automate many aspects of the container management, and provide the best solution for the scenario.

Quiz

You've been asked to write a program that uses Vision API to check for inappropriate content in photos that are uploaded to a Cloud Storage bucket. Any photos that are inappropriate should be deleted. What might be the simplest, cheapest way to deploy this program?

- A. Compute Engine
- B. GKE
- C. Cloud Functions
- D. App Engine



You've been asked to write a program that uses Vision API to check for inappropriate content in photos that are uploaded to a Cloud Storage bucket. Any photos that are inappropriate should be deleted. What might be the simplest, cheapest way to deploy this program?

- A. Compute Engine
- B. GKE
- C. Cloud Functions
- D. App Engine

Quiz

You've been asked to write a program that uses Vision API to check for inappropriate content in photos that are uploaded to a Cloud Storage bucket. Any photos that are inappropriate should be deleted. What might be the simplest, cheapest way to deploy this program?

- A. Compute Engine
- B. GKE
- C. Cloud Functions
- D. App Engine



C. This is the correct answer. The requirements for simplest and cheapest are met with Cloud Functions. Cloud Functions are for single purpose functions like image analysis. Cloud Functions also can be triggered by Cloud Storage events, so they provide seamless integration. The payment model based on number of requests, processing time of requests (measured in 100ms units), and then other resources consumed is the most suitable of all options offered above. There is a free tier too. Cloud Functions also provides automatic scaling, high availability, and fault tolerance.

Answers A, B, D could all be solutions but would require more development work and more expense for resources used, and as a result, do not meet the requirement of simplest, cheapest way of achieving the required functionality.

Review

Deploying Applications to Google Cloud



In this module we covered the various deployment services provided by Google. These include Compute Engine if you need complete control over your deployment environment; Google Kubernetes Engine if you want the flexibility, portability and automation that is provided by Kubernetes; and App Engine and Cloud Run if you want a completely managed platform as a service.

Again, each of these choices has advantages and disadvantages. Make sure you understand each one so that you can make an informed decision when deploying your services.

More resources

Migration to Google Cloud: Deploying your workloads

<https://cloud.google.com/solutions/migration-to-gcp-deploying-your-workloads>

Compute Engine

<https://cloud.google.com/compute/>

GKE

<https://cloud.google.com/kubernetes-engine/>

App Engine

<https://cloud.google.com/appengine/>



The links provide access to some useful resources on Google Cloud deployment platforms.

