



In this module, we cover virtual machine instances, or VMs.

VMs are the most common infrastructure component, and in GCP, they're provided by Compute Engine. A VM is similar, but not identical, to a hardware computer. VMs consists of a virtual CPU, some amount of memory, disk storage, and an IP address. Compute Engine is GCP's service to create VMs; it is very flexible and offers many options, including some that can't exist in physical hardware. For example, a micro VM shares a CPU with other virtual machines, so you can get a VM with less capacity at a lower cost. Another example of a function that can't exist in hardware is that some VMs offer burst capability, meaning that the virtual CPU will run above its rated capacity for a brief period, using the available shared physical CPU. The main VM options are CPUs, memory, disks, and networking.

Agenda

Compute Engine

Lab

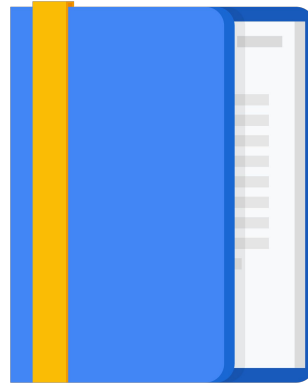
Compute Options (vCPU and
Memory)

Images

Disk Options

Common Compute Engine Actions

Lab



Now this is going to be a very robust module; there's a lot of detail to cover here with how virtual machines work on GCP. First we'll start with the basics of Compute Engine, followed by a quick little lab to get you more familiar with creating virtual machines.






Then, we'll look at the different CPU and memory options that enable you to create different configurations.

Next, we will look at images and the different disk options available with Compute Engine.

After that, we will discuss very common Compute Engine actions that you might encounter in your day-to-day job. This will be followed by an in-depth lab that explores many of the features and services covered in this module.

Let's get started with an overview of Compute Engine!

GCP compute and processing options

	 Compute Engine	 Kubernetes Engine	 App Engine Standard	 App Engine Flexible	 Cloud Functions
Language support	Any	Any	Python Node.js Go Java PHP	Python Node.js Go Java PHP Ruby .NET Custom Runtimes	Python Node.js Go
Usage model	IaaS	IaaS PaaS	PaaS	PaaS	Microservices Architecture
Scaling	Server Autoscaling	Cluster	Autoscaling managed servers		Serverless
Primary use case	General Workloads	Container Workloads	Scalable web applications Mobile backend applications		Lightweight Event Actions



As mentioned in the introduction to the course, there is a spectrum of different options in GCP for compute and processing. We will focus on the traditional virtual machine instances.

Now the difference is, Compute Engine gives you the utmost in flexibility: run whatever language you want—it's your virtual machine. This is purely an infrastructure as a service or IaaS model.

You have a VM and an operating system, and you can choose how to manage it and how to handle aspects, such as autoscaling, where you'll configure the rules about adding more virtual machines in specific situations. Autoscaling will be covered in a later course of this series.

The primary work case of Compute Engine is any general workload, especially an enterprise application that was designed to run on a server infrastructure. This makes Compute Engine very portable and easy to run in the cloud. Other services, like Google Kubernetes Engine, which consists of containerized workloads, may not be as easily transferable as what you're used to from on-premises.

Compute Engine

Infrastructure as a Service (IaaS)

Predefined or custom machine types:

- vCPUs (cores) and Memory (RAM)
- Persistent disks: HDD, SSD, and Local SSD
- Networking
- Linux or Windows

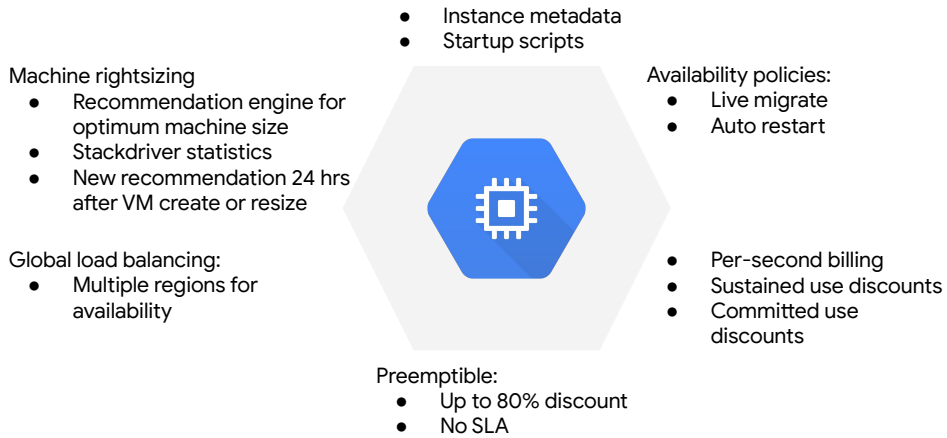


Compute Engine



So what is Compute Engine? At its heart, it's physical servers that you're used to, running inside the GCP environment, with a number of different configurations. Both predefined and custom machine types allow you to choose how much memory and how much CPU you want. You choose the type of disk you want, whether you want to just use standard hard drives, SSDs, local SSDs, or a mix. You can even configure the networking interfaces and run a combination of Linux and Windows machines.

Compute Engine features



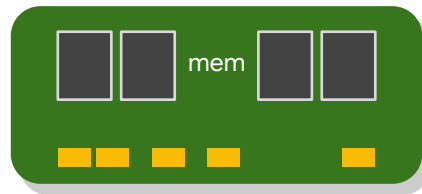
Several different features will be covered throughout this module, such as machine rightsizing, startup scripts, metadata, availability policies, and pricing and usage discounts.

Compute

Several machine types

- Network throughput scales 2 Gbps per vCPU (small exceptions)
- Theoretical max of 32 Gbps with 16 vCPU or 100 Gbps with T4 or V100 GPUs

A vCPU is equal to 1 hardware hyper-thread



Let's start by looking at the compute options.

Compute Engine provides several different machine types that we'll discuss later in this module. If those machines don't meet your needs, you can also customize your own machine.

Your choice of CPU will affect your network throughput. Specifically, your network will scale at 2 Gbits per second for each CPU core, except for instances with 2 or 4 CPUs which receive up to 10 Gbits per second of bandwidth.

As of this recording there is a theoretical maximum throughput of 32 Gbits per second for an instance with 16 or more CPUs and a 100 Gbits per second maximum throughput for specific instances that have T4 or V100 GPUs attached.

When you're migrating from an on-premises setup, you're used to physical cores, which have hyperthreading. On Compute Engine, each virtual CPU (or vCPU) is implemented as a single hardware hyper-thread on one of the available CPU Platforms.

For an up-to-date list of all the available CPU platforms, refer to the links section of this video [<https://cloud.google.com/compute/docs/cpu-platforms>]

Storage

Disks

- Standard, SSD, or Local SSD
- Standard and SSD PDs scale in performance for each GB of space allocated

Resize disks or migrate instances with no downtime



After you pick your compute options, you want to choose your disk.

You have three options: Standard, SSD, or local SSD. So basically, do you want the standard spinning hard disk drives (HDDs), or flash memory solid-state drives (SSDs)? Both of these options provide the same amount of capacity in terms of disk size when choosing a persistent disk. Therefore, the question really is about performance versus cost, because there's a different pricing structure.

Basically, SSDs are designed to give you a higher number of IOPS per dollar versus standard disks, which will give you a higher amount of capacity for your dollar.

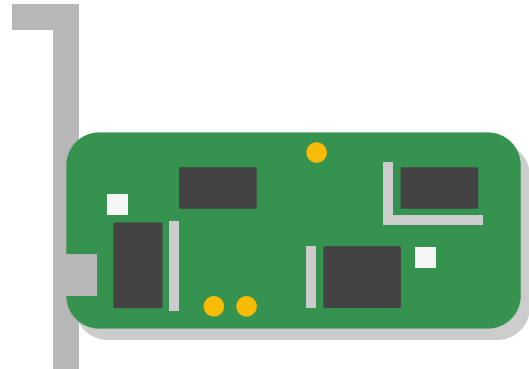
Local SSDs have even higher throughput and lower latency than SSD persistent disks, because they are attached to the physical hardware. However, the data that you store on local SSDs persists only until you stop or delete the instance. Typically, a local SSD is used as a swap disk, just like you would do if you want to create a ramdisk, but if you need more capacity, you can store those on a local SSD. You can create instances with up to eight separate 375-GB local SSD partitions for a total of 3 TB of local SSD space for each instance.

Standard and non-local SSD disks can be sized up to 64 TB for each instance. The performance of these disks scales with each GB of space allocated.

Networking

Robust networking features

- Default, custom networks
- Inbound/outbound firewall rules
 - IP based
 - Instance/group tags
- Regional HTTPS load balancing
- Network load balancing
 - Does not require pre-warming
- Global and multi-regional subnetworks



As for networking, we have already seen networking features applied to Compute Engine in the previous module's lab. We looked at the different types of networks and created firewall rules using IP addresses and network tags.

You'll also notice that you can do regional HTTPS load balancing and network load balancing. This doesn't require any pre-warming because a load balancer isn't a hardware device that needs to analyze your traffic. A load balancer is essentially a set of traffic engineering rules that are coming into the Google network, and VPC is applying your rules destined for your IP address subnet range. We'll learn more about load balancers in a later course of the Architecting with Google Compute Engine series.

Demo

Create a VM



Let me give you a quick walk through of the VM instance creation process and point out CPU, storage, and network options in the GCP Console.

[Demo]

That's how easy it is to configure the location, CPU, memory, storage, and network interface for a VM instance using the GCP Console. Let's get back to the slides to go over VM access and lifecycle.

VM access

Linux: SSH

- SSH from GCP Console or CloudShell via Cloud SDK
- SSH from computer or third-party client and generate key pair
- Requires firewall rule to allow tcp:22

Windows: RDP

- RDP clients
- Powershell terminal
- Requires setting the Windows password
- Requires firewall rule to allow tcp:3389



For accessing a VM, the creator of an instance has full root privileges on that instance.

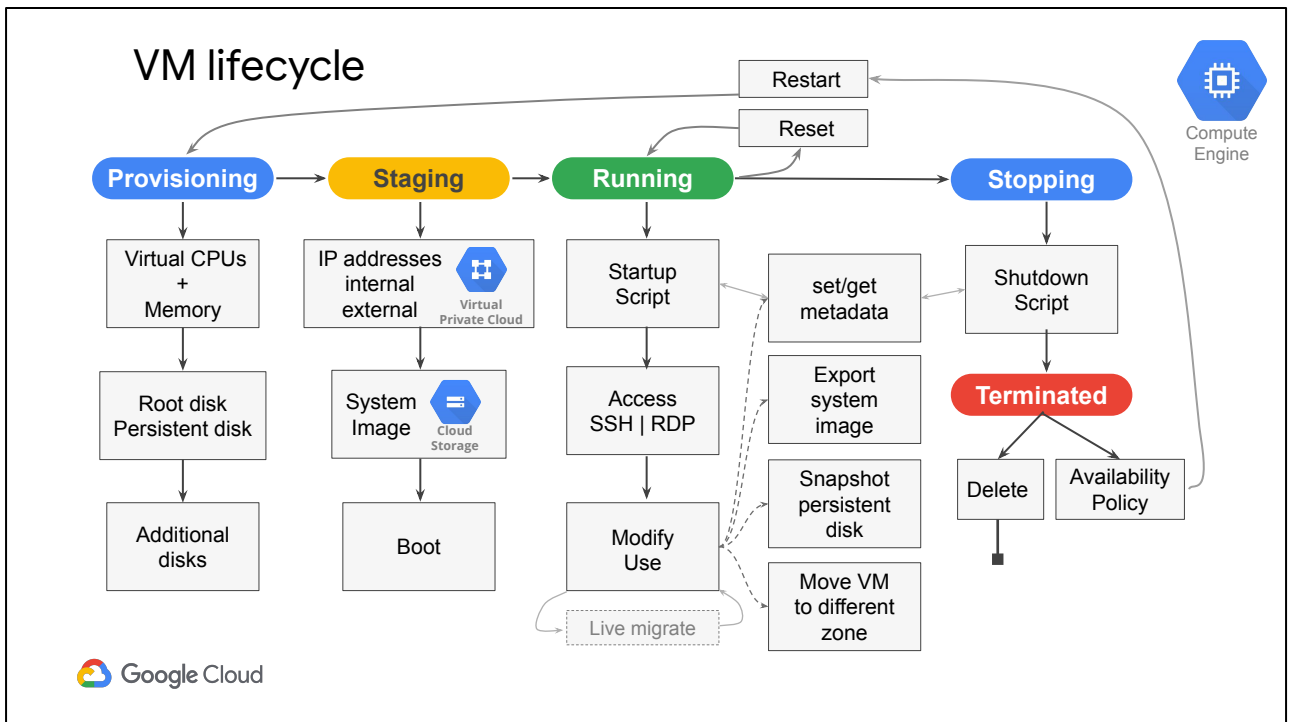
On a Linux instance, the creator has SSH capability and can use the GCP Console to grant SSH capability to other users.

On a Windows instance, the creator can use the GCP Console to generate a username and password. After that, anyone who knows the username and password can connect to the instance using a Remote Desktop Protocol, or RDP, client.

I listed the required firewall rules for both SSH and RDP here, but you don't need to define these if you are using the default network that we covered in the previous module.

For more information on SSH key management and creating passwords for Windows instances, refer to the links section of this video:

<https://cloud.google.com/compute/docs/instances/adding-removing-ssh-keys>
<https://cloud.google.com/compute/docs/instances/windows/creating-passwords-for-windows-instances>



The lifecycle of a VM is represented by different statuses. I will cover this lifecycle on a high level, but I recommend returning to this diagram as a reference.

When you define all the properties of an instance and click Create, the instance enters the provisioning state. Here the resources such as CPU, memory, and disks are being reserved for the instance, but the instance itself isn't running yet. Next, the instance moves to the staging state where resources have been acquired and the instance is prepared for launch. Specifically, in this state, Compute Engine is adding IP addresses, booting up the system image, and booting up the system.

After the instance starts running, it will go through pre-configured startup scripts and enable SSH or RDP access. Now, you can do several things while your instance is running. For example, you can live migrate your virtual machine to another host in the same zone instead of requiring your instance to be rebooted. This allows GCP to perform maintenance that is integral to keeping the infrastructure protected and reliable, without interrupting any of your VMs. While your instance is running, you can also move your VM to a different zone, take a snapshot of the VM's persistent disk, export the system image, or reconfigure metadata. We will explore some of these tasks in later labs.

Some actions require you to stop your virtual machine; for example, if you want to

upgrade your machine by adding more CPU. When the instance enters this state, it will go through pre-configured shutdown scripts and end in the terminated state. From this state, you can choose to either restart the instance, which would bring it back to its provisioning state, or delete it.

You also have the option to reset a VM, which is similar to pressing the reset button on your computer. This action wipes the memory contents of the machine and resets the virtual machine to its initial state. The instance remains in the running state through the reset.

Changing VM state from running

	methods	Shutdown Script time	state
reset	console, gcloud, API, OS	no	remains running
restart	console, gcloud, API, OS	no	terminated → running
reboot	OS: <code>sudo reboot</code>	~90 sec	running → running
stop	console, gcloud, API	~90 sec	running → terminated
shutdown	OS: <code>sudo shutdown</code>	~90 sec	running → terminated
delete	console, gcloud, API	~90 sec	running → N/A
<i>preemption</i>	<i>automatic</i>	~30 sec	N/A

"ACPI Power Off"



There are different ways you can change a VM state from running. Some methods involve the GCP Console and the gcloud command, while others are performed from the OS, such as for reboot and shutdown.

It's important to know that if you are restarting, rebooting, stopping, or even deleting an instance, the shutdown process will take about 90 sec. For a preemptible VM, if the instance does not stop after 30 seconds, Compute Engine sends an ACPI G3 Mechanical Off signal to the operating system. Remember that when writing shutdown scripts for preemptible VMs.

Availability policy: Automatic changes

Called "scheduling options" in SDK/API

Automatic restart

- Automatic VM restart due to crash or maintenance event
 - Not preemption or a user-initiated terminate

On host maintenance

- Determines whether host is live-migrated or terminated due to a maintenance event. Live migration is the default.

Live migration

- During maintenance event, VM is migrated to different hardware without interruption.
- Metadata indicates occurrence of live migration.



As I mentioned previously, Compute Engine can live migrate your virtual machine to another host due to a maintenance event to prevent your applications from experiencing disruptions. A VM's availability policy determines how the instance behaves in such an event.

The default maintenance behavior for instances is to live migrate, but you can change the behavior to terminate your instance during maintenance events instead. If your VM is terminated due to a crash or other maintenance event, your instance automatically restarts by default, but this can also be changed.

These availability policies can be configured both during the instance creation and while an instance is running by configuring the Automatic restart and On host maintenance options.

For more information on live migration, refer to the links section of this video:

<https://cloud.google.com/compute/docs/instances/live-migration>

Stopped (Terminated) VM

No charge for stopped VM

- Charged for attached disks and IPs

Actions

- Change the machine type.
- Add or removed attached disks; change auto-delete settings.
- Modify instance tags.
- Modify custom VM or project-wide metadata.
- Remove or set a new static IP.
- Modify VM availability policy.
- Can't change the image of a stopped VM.



When a VM is terminated, you do not pay for memory and CPU resources. However, you are charged for any attached disks and reserved IP addresses. In the terminated state, you can perform any of the actions listed here, such as changing the machine type, but you cannot change the image of a stopped VM.

Also, not all of the actions listed here require you to stop a virtual machine. For example, VM availability policies can be changed while the VM is running, as discussed previously.

Lab

Creating Virtual Machines



Let's take some of the Compute Engine concepts that we just discussed and apply them in a lab.

In this lab, you explore virtual machine instance options by creating several standard VMs and a custom VM. You also connect to those VMs using both SSH for Linux machines and RDP for Windows machines.

Lab review

Creating Virtual Machines



In this lab, you created several virtual machine instances of different types with different characteristics. Specifically, you created a small utility VM for administration purposes, a Windows VM, and a custom Linux VM. You also accessed both the Windows and Linux VMs and deleted all your created VMs.

In general, start with smaller VMs when you're prototyping solutions to keep the costs down. When you are ready for production, trade up to larger VMs based on capacity. If you're building in redundancy for availability, remember to allocate excess capacity to meet performance requirements. Finally, consider using custom VMs when your application's requirements fit between the features of the standard types.

You can stay for a lab walkthrough, but remember that GCP's user interface can change, so your environment might look slightly different.

Agenda

Compute Engine

Lab

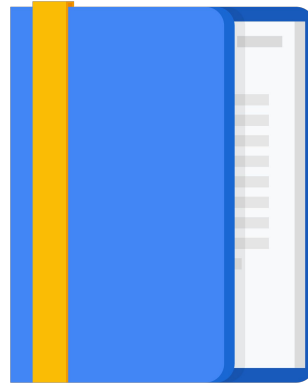
Compute Options (vCPU and
Memory)

Images

Disk Options

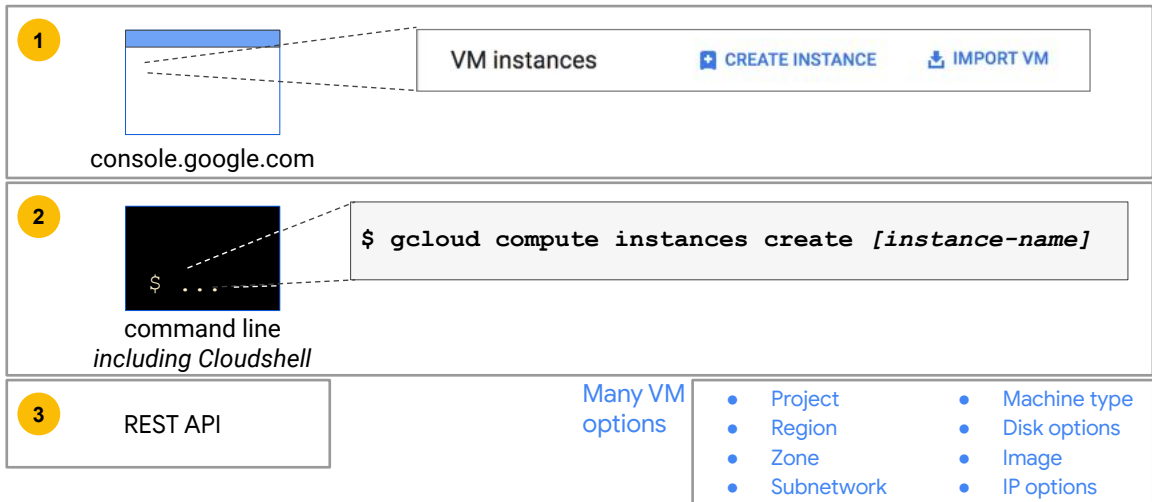
Common Compute Engine Actions

Lab



Now that you have completed the lab, let's dive deeper into the compute options that are available to you in GCP, by focusing on CPU and memory.

Creating a VM



You have three options for creating and configuring a VM. You can use the GCP Console as you did in the previous lab, the Cloud Shell command line, or the RESTful API. If you'd like to automate and process very complex configurations, you might want to programmatically configure these through the RESTful API by defining all the different options for your environment.

If you plan on using the command line or RESTful API, I recommend that you first configure the instance through the GCP Console and then ask Compute Engine for the equivalent REST request or command line, as I showed you in my demo earlier. This way you avoid any typos and get dropdown lists of all the available CPU and memory options.

Speaking of CPU and memory options, let's look at the different machine types that are currently available.

Machine types

Predefined machine types: Ratio of GB of memory per vCPU

- Standard
- High-memory
- High-CPU
- Memory-optimized
- Compute-optimized
- Shared-core

Custom machine types:

- You specify the amount of memory and number of vCPUs.



A machine type specifies a particular collection of virtualized hardware resources available to a VM instance, including the system memory size, vCPU count, and maximum persistent disk capability. GCP offers several machine types that can be grouped into 2 categories:

- Predefined machine types: These have a fixed collection of resources, are managed by Compute Engine and are available in multiple different classes. Each class has a predefined ratio of GB of memory per vCPU. These are the:
 - Standard machine types
 - High-memory machine types
 - High-CPU machine types
 - Memory-optimized machine types
 - Compute-optimized machine types
 - Shared-core machine types
- Custom machine types: These let you specify the number of vCPUs and the amount of memory for your instance.

Let's explore each of these machine types, but remember that these machine types and their available options can change.

Standard machine types

3.75 GB of memory
1 vCPU

Machine name	vCPUs	Memory (GB)	Max # PD	Max total PD size
n1-standard-1	1	3.75	128	64 TB
n1-standard-2	2	7.50		
n1-standard-4	4	15		
n1-standard-8	8	30		
n1-standard-16	16	60		
n1-standard-32	32	120		
n1-standard-64	64	240		
n1-standard-96	96	360		



Standard machine types are suitable for tasks that have a balance of CPU and memory needs. Standard machine types have 3.75 GB of memory per vCPU. The vCPU configurations come in different intervals from 1 vCPU all the way to 96 vCPUs, as shown on this table. Each of these machines supports a maximum of 128 persistent disks with a total persistent disk size of 64 TB, which is also the case for the High-memory, High-CPU, Memory-optimized, and Compute-optimized machine types.

High-memory machine types

6.5 GB of memory
1 vCPU

Machine name	vCPUs	Memory (GB)	Max # PD	Max total PD size
n1-highmem-2	2	13	128	64 TB
n1-highmem-4	4	26		
n1-highmem-8	8	52		
n1-highmem-16	16	104		
n1-highmem-32	32	208		
n1-highmem-64	64	416		
n1-highmem-96	96	624		



High-memory machine types are ideal for tasks that require more memory relative to vCPUs. High-memory machine types have 6.50 GB of system memory per vCPU. Similarly to the standard machine types, the vCPU configurations come in different intervals from 2 vCPUs all the way to 96 vCPUs, as shown on this table.

High-CPU machine types

0.9 GB of memory

1 vCPU

Machine name	vCPUs	Memory (GB)	Max # PD	Max total PD size
n1-highcpu-2	2	1.80	128	64 TB
n1-highcpu-4	4	3.60		
n1-highcpu-8	8	7.20		
n1-highcpu-16	16	14.4		
n1-highcpu-32	32	28.8		
n1-highcpu-64	64	57.6		
n1-highcpu-96	96	86.4		



High-CPU machine types are ideal for tasks that require more vCPUs relative to memory. High-CPU machine types have 0.90 GB of memory per vCPU.

Memory-optimized machine types

<div>>14 GB of memory</div> <div>1 vCPU</div>	Machine name	vCPUs	Memory (GB)	Max # PD	Max total PD size
	n1-ultramem-40	40	961	128	64 TB
	n1-ultramem-80	80	1922		
	n1-megamem-96	96	1433.6		
	n1-ultramem-160	160	3844		



Memory-optimized machine types are ideal for tasks that require intensive use of memory, with higher memory to vCPU ratios than high-memory machine types. These machines types are perfectly suited for in-memory databases and in-memory analytics, such as SAP HANA and business warehousing workloads, genomics analysis, and SQL analysis services. Memory-optimized machine types have more than 14 GB of memory per vCPU.

These machine come in 4 configurations as shown in this table, with only the n1-megamem-96 supporting a local SSD as of this recording.

Compute-optimized machine types

Highest performance per vCPU (3.8Ghz sustained all-core turbo)

Machine name	vCPUs	Memory (GB)	Max # PD	Max total PD size
c2-standard-4	4	16	128	64 TB
c2-standard-8	8	32		
c2-standard-16	16	64		
c2-standard-30	30	120		
c2-standard-60	60	240		



Compute-optimized machine types are ideal for compute-intensive workloads. These machine types offer the highest performance per core on Compute Engine. Built on the latest-generation Intel Scalable Processors (the Cascade Lake), C2 machine types offer up to 3.8 Ghz sustained all-core turbo and provide full transparency into the architecture of the underlying server platforms, enabling advanced performance tuning. C2 machine types offer much more computing power, run on a newer platform, and are generally more robust for compute-intensive workloads than the N1 high-CPU machine types.

Shared-core machine types

Machine name	vCPUs	Memory (GB)	Max # PD	Max total PD size
f1-micro	0.2	0.60	16	3 TB
g1-small	0.5	1.70		



Shared-core machine types provide one vCPU that is allowed to run for a portion of the time on a single hardware hyper-thread on the host CPU running your instance. Shared-core instances can be more cost-effective for running small, non-resource-intensive applications than other machine types. There are only two shared-core machine types to choose from:

- f1-micro
- g1-small

f1-micro machine types offer bursting capabilities that allow instances to use additional physical CPU for short periods of time. Bursting happens automatically when your instance requires more physical CPU than originally allocated. During these spikes, your instance will opportunistically take advantage of available physical CPU in bursts. Note that bursts are not permanent and are only possible periodically.

For up-to-date information about all of these machine types, see the links section of this video: [\[https://cloud.google.com/compute/docs/machine-types\]](https://cloud.google.com/compute/docs/machine-types)

Creating custom machine types

When to select custom:

- Requirements fit between the predefined types
- Need more memory or more CPU

Customize the amount of memory and vCPU for your machine:

- Either 1 vCPU or even number of vCPU
- 0.9 GB per vCPU, up to 6.5 GB per vCPU (default)
- Total memory must be multiple of 256 MB

Machine type
Customize to select cores, memory and GPUs.

[Basic view](#)

Cores

1 vCPU 1 - 96

Memory

3.75 GB 1 - 6.5

☐ Extend memory ?

CPU platform ?

Automatic

GPUs

The number of GPU dies is linked to the number of CPU cores and memory selected for this instance. For the current configuration, you can select no fewer than 1 GPU die of this type. [Learn more](#)

Number of GPUs **GPU type**

None NVIDIA Tesla K80



If none of the predefined machine types match your needs, you can independently specify the number of vCPUs and the amount of memory for your instance. Custom machine types are ideal for the following scenarios:

- When you have workloads that are not a good fit for the predefined machine types that are available to you.
- Or when you have workloads that require more processing power or more memory, but don't need all of the upgrades that are provided by the next larger predefined machine type.

It costs slightly more to use a custom machine type than an equivalent predefined machine type, and there are still some limitations in the amount of memory and vCPUs you can select:

- Only machine types with 1 vCPU or an even number of vCPUs can be created.
- Memory must be between 0.9 GB and 6.5 GB per vCPU (by default).
- The total memory of the instance must be a multiple of 256 MB.

By default, a custom machine can have up to 6.5 GB of memory per vCPU. However, this might not be enough memory for your workload. At an additional cost, you can get more memory per vCPU beyond the 6.5 GB limit. This is referred to as *extended memory*, and you can learn more about this in the links section of this video

[\[https://cloud.google.com/compute/docs/instances/creating-instance-with-custom-machine-type#extendedmemory\]](https://cloud.google.com/compute/docs/instances/creating-instance-with-custom-machine-type#extendedmemory)

Choosing region and zone



The first thing you want to consider when choosing a region and zone is the geographical location in which you want to run your resources. This map shows the current and planned GCP regions and number of zones. For up-to-date information on the available regions and zones, see the documentation linked for this video: [\[https://cloud.google.com/compute/docs/regions-zones/#available\]](https://cloud.google.com/compute/docs/regions-zones/#available)

Each zone supports a combination of Ivy Bridge, Sandy Bridge, Haswell, Broadwell, and Skylake platforms. When you create an instance in the zone, your instance will use the default processor supported in that zone. For example, if you create an instance in the us-central1-a zone, your instance will use a Sandy Bridge processor.

Pricing

- Per-second billing, with minimum of 1 minute
 - vCPUs, GPUs, and GB of memory
- Resource-based pricing:
 - Each vCPU and each GB of memory is billed separately
- Discounts:
 - Sustained use
 - Committed use
 - Preemptible VM instances
- Recommendation Engine
 - Notifies you of underutilized instances
- Free usage limits



GCP offers a variety of different options to keep the prices low for Compute Engine resources:

All vCPUs, GPUs, and GB of memory are charged a minimum of 1 minute. For example, if you run your virtual machine for 30 seconds, you will be billed for 1 minute of usage. After 1 minute, instances are charged in 1-second increments.

Compute Engine uses a resource-based pricing model, where each vCPU and each GB of memory on Compute Engine is billed separately rather than as part of a single machine type. You still create instances using predefined machine types, but your bill reports them as individual vCPUs and memory used.

There are several discounts available but the discount types cannot be combined:

- Resource-based pricing allows Compute Engine to apply sustained use discounts to all of your predefined machine type usage in a region collectively rather than to individual machine types.
- If your workload is stable and predictable, you can purchase a specific amount of vCPUs and memory for a discount off of normal prices in return for committing to a usage term of 1 year or 3 years. The discount is up to 57% for most machine types or custom machine types. The discount is up to 70% for memory-optimized machine types.

- A preemptible VM is an instance that you can create and run at a much lower price than normal instances. However, Compute Engine might terminate (or preempt) these instances if it requires access to those resources for other tasks. Preemptible instances are excess Compute Engine capacity so their availability varies with usage.

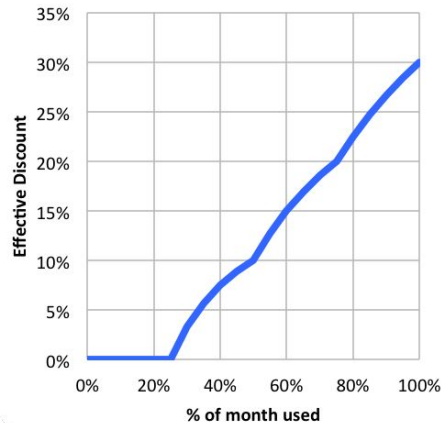
The ability to customize the amount of memory and CPU through custom machine types allows for further pricing customization. Speaking of sizing your machine, Compute Engine provides VM sizing recommendations to help you optimize the resource use of your virtual machine instances. When you create a new instance, recommendations for the new instance will appear 24 hours after the instance has been created.

Compute Engine also has Free Usage Limits. For the exact terms, please refer to the links section of this video:

[\[https://cloud.google.com/free/docs/gcp-free-tier#always-free-usage-limits\]](https://cloud.google.com/free/docs/gcp-free-tier#always-free-usage-limits)

Sustained use discounts

Usage Level (% of month)	% at which incremental is charged
0% - 25%	100% of base rate
25% - 50%	80% of base rate
50% - 75%	60% of base rate
75% - 100%	40% of base rate



Up to 30% net discount for instances that run the entire month



Sustained use discounts are automatic discounts that you get for running specific Compute Engine resources (vCPUs, memory, GPU devices) for a significant portion of the billing month. For example, when you run one of these resources for more than 25% of a month, Compute Engine automatically gives you a discount for every incremental minute you use for that instance. The discount increases with usage, and you can get up to a 30% net discount for instances that run the entire month.

This table shown on this slide describes the discount you get at each usage level of a VM instance. To take advantage of the full 30% discount, create your VM instances on the first day of the month, because discounts reset at the beginning of each month.

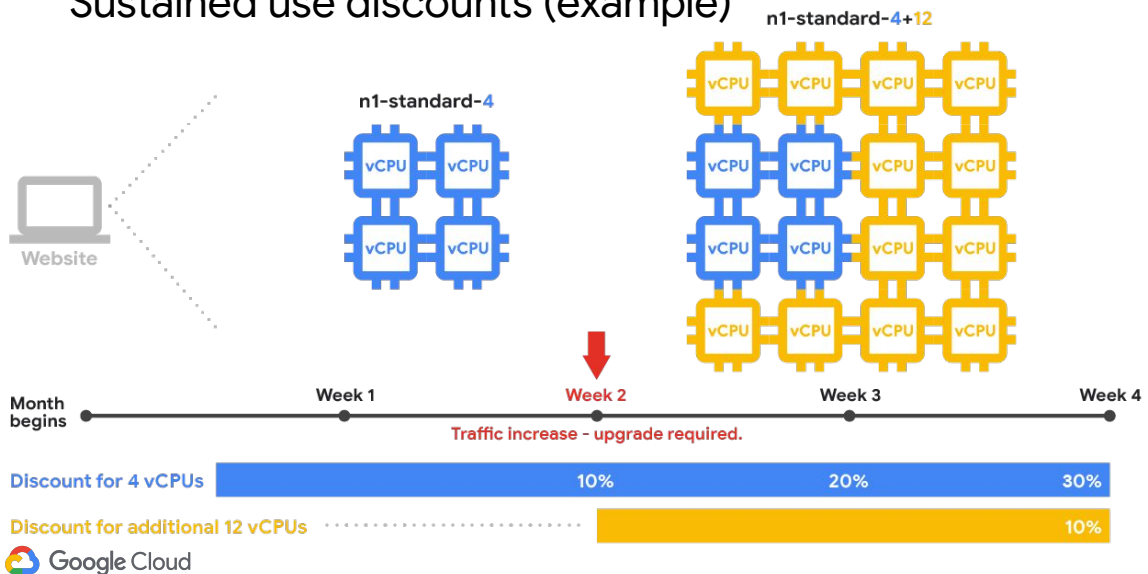
The graph on this slide demonstrates how your effective discount increases with use. For example, if you use a virtual machine for 50% of the month, you get an effective discount of 10%. If you use it for 75% of the month, you get an effective discount of 20%. If you use it for 100% of the month, you get an effective discount of 30%. You can also use the GCP Pricing Calculator to estimate your sustained use discount for any arbitrary workload. For the calculator, see the links section of this video:

[\[https://cloud.google.com/products/calculator/\]](https://cloud.google.com/products/calculator/)

Compute Engine calculates sustained use discounts based on vCPU and memory usage across each region and separately for each of the following categories:

- Predefined machine types
- Custom machine type

Sustained use discounts (example)



Let's go through an example where you have two instances that are in the same region but have different machine types and run at different times of the month. Compute Engine breaks down the number of vCPUs and amount of memory used across all instances that use predefined machine types and combines the resources to qualify for the largest sustained usage discounts possible. As shown on this slide, you run the following two instances in the us-central1 region during a month:

- For the first half of the month, you run an n1-standard-4 instance with 4 vCPUs and 15 GB of memory
- For the second half of the month, you run a larger n1-standard-16 instance with 16 vCPUs and 60 GB of memory

In this scenario, Compute Engine reorganizes these machine types into individual vCPUs and memory resources and combines their usage to create the following resources, as shown on the bottom:

- 4 vCPUs and 15 GB of memory for a full month
- 12 vCPUs and 45 GB of memory for half of the month

Preemptible

- Lower price for interruptible service (up to 80%)
- VM might be terminated at any time
 - No charge if terminated in the first 10 minutes
 - 24 hours max
 - 30-second terminate warning, but not guaranteed
 - *Time for a shutdown script*
- No live migrate; no auto restart
- You can request that CPU quota for a region be split between regular and preemption
 - Default: preemptible VMs count against region CPU quota



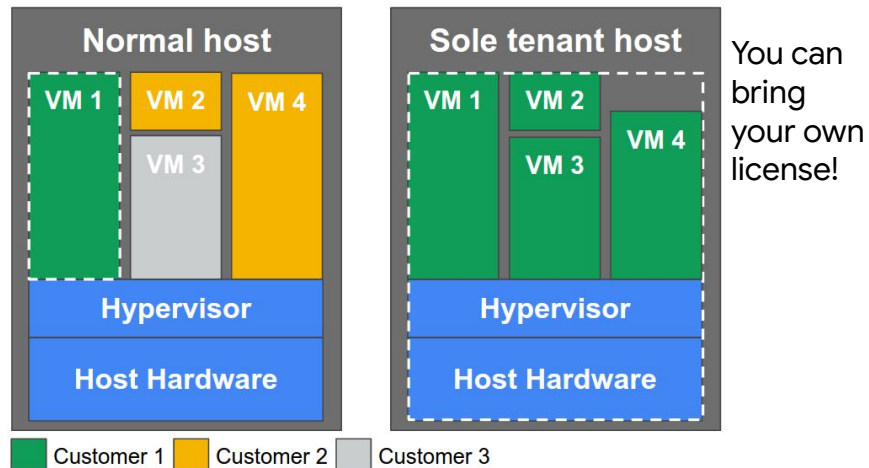
As I mentioned earlier, a preemptible VM is an instance that you can create and run at a much lower price than normal instances. See whether you can make your application function completely on preemptible VMs, because an 80-percent discount is a significant investment in your application. Now, just to reiterate, these VMs might be preempted at any time, and there is no charge if that happens within the first 10 minutes. Also, preemptible VMs are only going to live for up to 24 hours, and you only get a 30-second notification before the machine is preempted.

It's also worth noting that there are no live migrations nor automatic restarts in preemptible VMs, but something that we will highlight is that you can actually create monitoring and load balancers that can start up new preemptible VMs in case of a failure. In other words, there are external ways to keep restarting preemptible VMs if you need to.

One major use case for preemptible VMs is running a batch processing job. If some of those instances terminate during processing, the job slows but does not completely stop. Therefore, preemptible instances complete your batch processing tasks without placing additional workload on your existing instances, and without requiring you to pay full price for additional normal instances.

[\[https://cloud.google.com/compute/docs/instances/preemptible#what_is_a_preemptible_instance\]](https://cloud.google.com/compute/docs/instances/preemptible#what_is_a_preemptible_instance)

Sole-tenant nodes physically isolate workloads



If you have workloads that require physical isolation from other workloads or virtual machines in order to meet compliance requirements, you want to consider sole-tenant nodes.

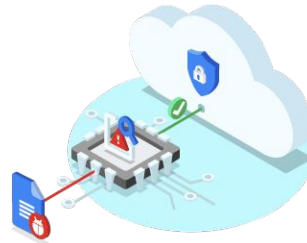
A sole-tenant node is a physical Compute Engine server that is dedicated to hosting VM instances only for your specific project. Use sole-tenant nodes to keep your instances physically separated from instances in other projects, or to group your instances together on the same host hardware, for example if you have a payment processing workload that needs to be isolated to meet compliance requirements.

The diagram on the left shows a normal host with multiple VM instances from multiple customers. A sole tenant node as shown on the right also has multiple VM instances, but they all belong to the same project. As of this recording, the only available node type can accommodate VM instances up to 96 vCPUs and 624 GB of memory. You can also fill the node with multiple smaller VM instances of various sizes, including custom machine types and instances with extended memory. Also, if you have existing operating system licenses, you can bring them to Compute Engine using sole-tenant nodes while minimizing physical core usage with the in-place restart feature

To learn how to create nodes and place your instances on those nodes, see the links section of this video: [\[https://cloud.google.com/compute/docs/nodes/create-nodes\]](https://cloud.google.com/compute/docs/nodes/create-nodes)

Shielded VMs offer verifiable integrity

- Secure Boot
- Virtual trusted platform module (vTPM)
- Integrity monitoring



Requires shielded image!



Another compute option is to create shielded VMs. Shielded VMs offer verifiable integrity of your VM instances, so you can be confident that your instances haven't been compromised by boot- or kernel-level malware or rootkits. Shielded VM's verifiable integrity is achieved through the use of Secure Boot, virtual trusted platform module or vTPM-enabled Measured Boot, and integrity monitoring.

Shielded VMs is the first offering in the Shielded Cloud initiative. The Shielded Cloud initiative is meant to provide an even more secure foundation for all of GCP by providing verifiable integrity and offering features, like vTPM shielding or sealing, that help prevent data exfiltration.

In order to use these shielded VM features, you need to select a shielded image. We'll learn about images in the next section.

Agenda

Compute Engine

Lab

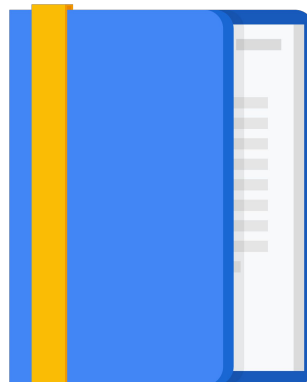
Compute Options (vCPU and
Memory)

Images

Disk Options

Common Compute Engine Actions

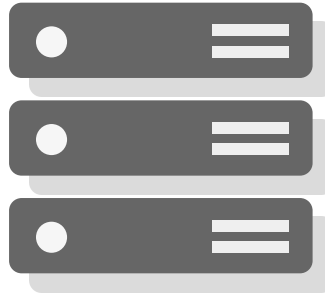
Lab



Next, let's focus on images.

What's in an image?

- Boot loader
- Operating system
- File system structure
- Software
- Customizations



When creating a virtual machine, you can choose the boot disk image. This image includes the boot loader, the operating system, the file system structure, any pre-configured software, and any other customizations.

Images

- Public base images
 - Google, third-party vendors, and community; Premium images (p)
 - Linux
 - CentOS, CoreOS, Debian, RHEL(p), SUSE(p), Ubuntu, openSUSE, and FreeBSD
 - Windows
 - Windows Server 2019(p), 2016(p), 2012-r2(p)
 - SQL Server pre-installed on Windows(p)
- Custom images
 - Create new image from VM: pre-configured and installed SW
 - Import from on-prem, workstation, or another cloud
 - Management features: image sharing, image family, deprecation



You can select either a public or custom image.

As you saw in the previous lab, you can choose from both Linux and Windows images. Some of these images are premium images, as indicated in parentheses with a p. These images will have per-second charges after a 1-minute minimum, with the exception of SQL Server images, which are charged per minute after a 10-minute minimum. Premium image prices vary with the machine type. However, these prices are global and do not vary by region or zone.

You can also use custom images. For example, you can create and use a custom image by pre-installing software that's been authorized for your particular organization.

You also have the option of importing images from your own premises or workstation, or from another cloud provider. This is a no-cost service that is as simple as installing an agent, and I highly recommend that you look at it. You can also share custom images with anybody in your project or among other projects, too.

Agenda

Compute Engine

Lab

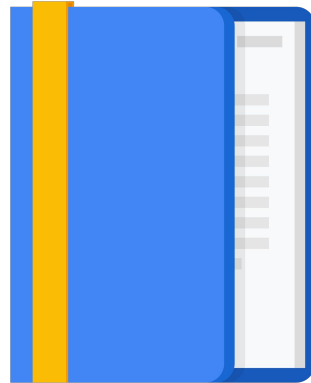
Compute Options (vCPU and
Memory)

Images

Disk Options

Common Compute Engine Actions

Lab



At this point you've chosen an operating system, but that operating system is going to be included as part of some kind of disk. So let's look at the disk options.

Boot disk

- VM comes with a single root persistent disk.
- Image is loaded onto root disk during first boot:
 - Bootable: you can attach to a VM and boot from it.
 - Durable: can survive VM terminate.
- Some OS images are customized for Compute Engine.
- Can survive VM deletion if “Delete boot disk when instance is deleted” is disabled.



Every single VM comes with a single root persistent disk, because you're choosing a base image to have that loaded on.

This image is bootable in that you can attach it to a VM and boot from it, and it is durable in that it can survive if the VM terminates. To have a boot disk survive a VM deletion, you need to disable the “Delete boot disk when instance is deleted” option in the instance’s properties.

As I discussed earlier, there are different types of disks. Let’s explore these in more detail.

Persistent disks

Network storage appearing as a block device

- Attached to a VM through the network interface
- Durable storage: *can* survive VM terminate
- Bootable: you can attach to a VM and boot from it
- Snapshots: incremental backups
- Performance: Scales with size

Features

- HDD (magnetic) or SSD (faster, solid-state) options
- Disk resizing: even running and attached!
- Can be attached in read-only mode to multiple VMs
- Encryption keys:
 - Google-managed
 - Customer-managed
 - Customer-supplied



The first disk that we create is what we call a persistent disk. That means it's going to be attached to the VM through the network interface. Even though it's persistent, it's not physically attached to the machine. This separation of disk and compute allows the disk to survive if the VM terminates. You can also perform snapshots of these disks, which are incremental backups that we'll discuss later.

The choice between HDD and SSD disks comes down to cost and performance. To learn more about disk performance and how it scales with disk size, see the links section of this video: [<https://cloud.google.com/compute/docs/disks/performance>]

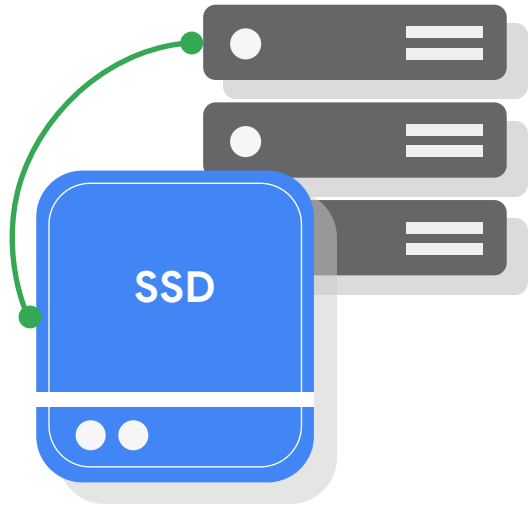
Another cool feature of persistent disks is that you can dynamically resize them, even while they are running and attached to a VM.

You can also attach a disk in read-only mode to multiple VMs. This allows you to share static data between multiple instances, which is cheaper than replicating your data to unique disks for individual instances.

By default, Compute Engine encrypts all data at rest. GCP handles and manages this encryption for you without any additional actions on your part. However, if you wanted to control and manage this encryption yourself, you can either use Cloud Key Management Service to create and manage key encryption keys (which is known as customer-managed encryption keys) or create and manage your own key encryption keys (known as customer-supplied encryption keys).

Local SSD disks are physically attached to a VM

- More IOPS, lower latency, and higher throughput than persistent disk
- 375-GB disk up to eight, total of 3 TB
- Data survives a reset, but not a VM stop or terminate
- VM-specific: cannot be reattached to a different VM



Now, local SSDs are different from persistent disks in that they are physically attached to the virtual machine. Therefore, these disks are ephemeral but provide very high IOPS. For up-to-date numbers I recommend referring to the documentation [\[https://cloud.google.com/compute/docs/disks/performance\]](https://cloud.google.com/compute/docs/disks/performance).

Currently, you can attach up to 8 local SSD disks with 375 GB each, resulting in a total of 3 TB.

Data on these disks will survive a reset but not a VM stop or terminate, because these disks can't be reattached to a different VM.

RAM disk

- `tmpfs`
- Faster than local disk, slower than memory
 - Use when your application expects a file system structure and cannot directly store its data in memory
 - Fast scratch disk, or fast cache
- Very volatile; erase on stop or restart
- May need a larger machine type if RAM was sized for the application
- Consider using a persistent disk to back up RAM disk data



You also have the option of using a RAM disk.

You can simply use `tmpfs` if you want to store data in memory. This will be the fastest type of performance available if you need small data structures. I recommend a high-memory virtual machine if you need to take advantage of such features, along with a persistent disk to back up the RAM disk data.

Summary of disk options

	Persistent disk HDD	Persistent disk SSD	Local SSD disk	RAM disk
Data redundancy	Yes	Yes	No	No
Encryption at rest	Yes	Yes	Yes	N/A
Snapshotting	Yes	Yes	No	No
Bootable	Yes	Yes	No	Not
Use case	General, bulk file storage	Very random IOPS	High IOPS and low latency	low latency and risk of data loss



In summary, you have several different disk options. Persistent disks can be rebooted and snapshotted, but local SSDs and RAM disks are ephemeral.

I recommend choosing a persistent HDD disk when you don't need performance but just need capacity. If you have high performance needs, start looking at the SSD options. The persistent disks offer data redundancy because the data on each persistent disk is distributed across several physical disks.

Local SSDs provide even higher performance, but without the data redundancy.

Finally, RAM disks are very volatile but they provide the highest performance.

Maximum persistent disks

Machine Type	Disk number limit
Shared-core	16
Standard	128
High-memory	
High-CPU	
Memory-optimized	
Compute-optimized	



Now, just as there is a limit on how many Local SSDs you can attach to a VM, there is also a limit on how many persistent disks you can attach to a VM. As illustrated in this table, this limit depends on the machine type. For the Shared-core machine type, you can attach up to 16 disks. For the Standard, High Memory, High-CPU, Memory-optimized, and Compute-optimized machine types, you can attach up to 128 disks. So you can create massive amounts of capacity for a single host.

Now remember that little nuance when I told you about how throughput is limited by the number of cores that you have? That throughput also shares the same bandwidth with Disk IO. So if you plan on having a large amount of Disk IO throughput, it will also compete with any network egress or ingress throughput. So remember that, especially if you will be increasing the number of drives attached to a virtual machine.

Persistent disk management differences

Cloud Persistent Disk

- Single file system is best
- Resize (grow) disks
- Resize file system
- Built-in snapshot service
- Automatic encryption



Computer Hardware Disk

- Partitioning
- Repartition disk
- Reformat
- Redundant disk arrays
- Subvolume management and snapshots
- Encrypt files before write to disk



There are many differences between a physical hard disk in a computer and a persistent disk, which is essentially a virtual networked device. First of all, if you remember with normal computer hardware disks, you have to partition them. Essentially, you have a drive and you're carving up a section for the operating system to get its own capacity. If you want to grow it, you have to repartition, and if you want to make changes you might even have to reformat. If you want redundancy, you might create a redundant disk array, and if you want encryption, you need to encrypt files before writing them to the disk.

With cloud persistent disks, things are very different because all that management is handled for you on the backend. You can simply grow disks and resize the file system because disks are virtual networked devices. Redundancy and snapshot services are built in and disks are automatically encrypted. You can even use your own keys, and that will ensure that no party can get to the data except you.

Agenda

Compute Engine

Lab

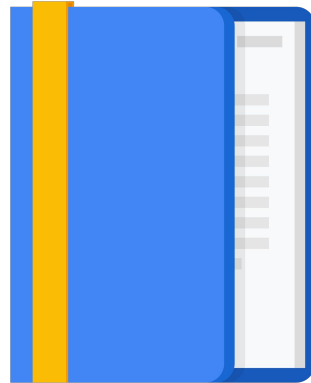
Compute Options (vCPU and
Memory)

Images

Disk Options

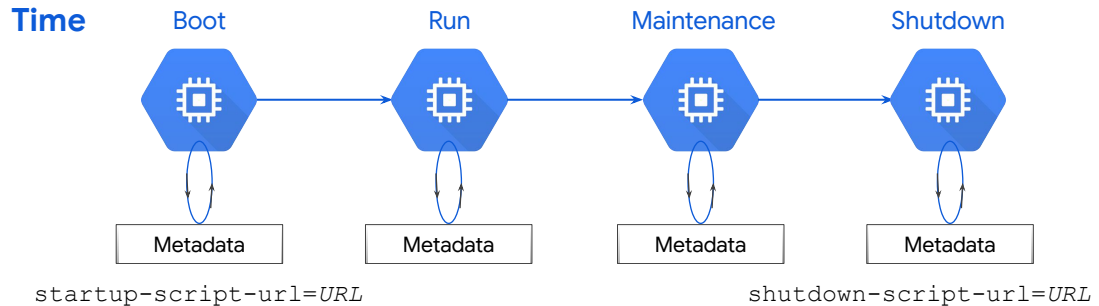
Common Compute Engine Actions

Lab



Now that we have covered all the different compute, image, and disk options, let's look at some common actions that you can perform with Compute Engine.

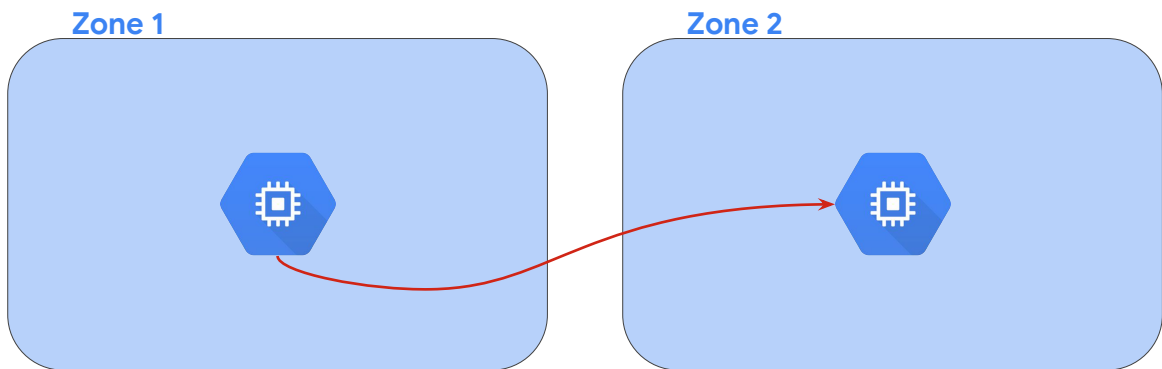
Metadata and scripts



Every VM instance stores its metadata on a metadata server. The metadata server is particularly useful in combination with startup and shutdown scripts, because you can use the metadata server to programmatically get unique information about an instance, without additional authorization. For example, you can write a startup script that gets the metadata key/value pair for an instance's external IP address and use that IP address in your script to set up a database. Because the default metadata keys are the same on every instance, you can reuse your script without having to update it for each instance. This helps you create less brittle code for your applications.

Storing and retrieving instance metadata is a very common Compute Engine action. I recommend storing the startup and shutdown scripts in Cloud Storage, as you will explore in the upcoming lab of this module.

Move an instance to a new zone



`gcloud compute instances move`



Another common action is to move an instance to a new zone. For example, you might do so for geographical reasons or because a zone is being deprecated.

Move an instance to a new zone

- Automated process (moving within region):
 - `gcloud compute instances move`
 - Update references to VM; not automatic
- Manual process (moving between regions):
 - Snapshot all persistent disks on the source VM.
 - Create new persistent disks in destination zone restored from snapshots.
 - Create new VM in the destination zone and attach new persistent disks.
 - Assign static IP to new VM.
 - Update references to VM.
 - Delete the snapshots, original disks, and original VM.

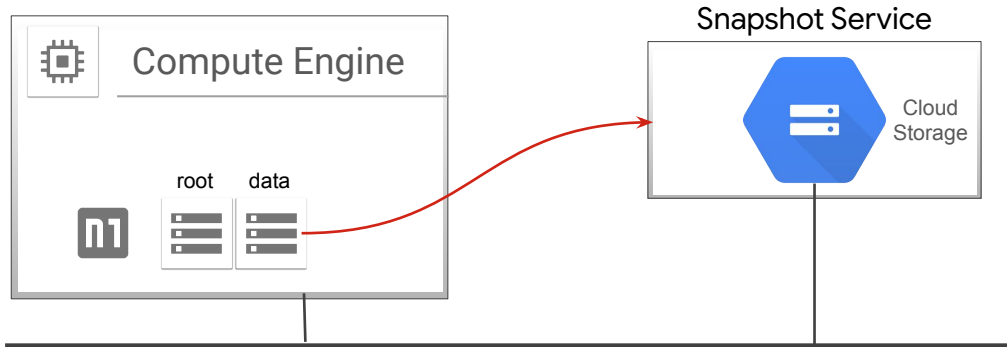


If you move your instance within the same region, you can automate the move by using the `gcloud compute instances move` command.

If you move your instance to a different region, you need to manually do so by following the process outlined here. This involves making a snapshot of all persistent disks and creating new disks in the destination zone from that snapshot. Next, you create the new VM in the destination zone and attach the new persistent disks, assign a static IP, and update any references to the VM. Finally, you delete the original VM, its disks, and the snapshot.

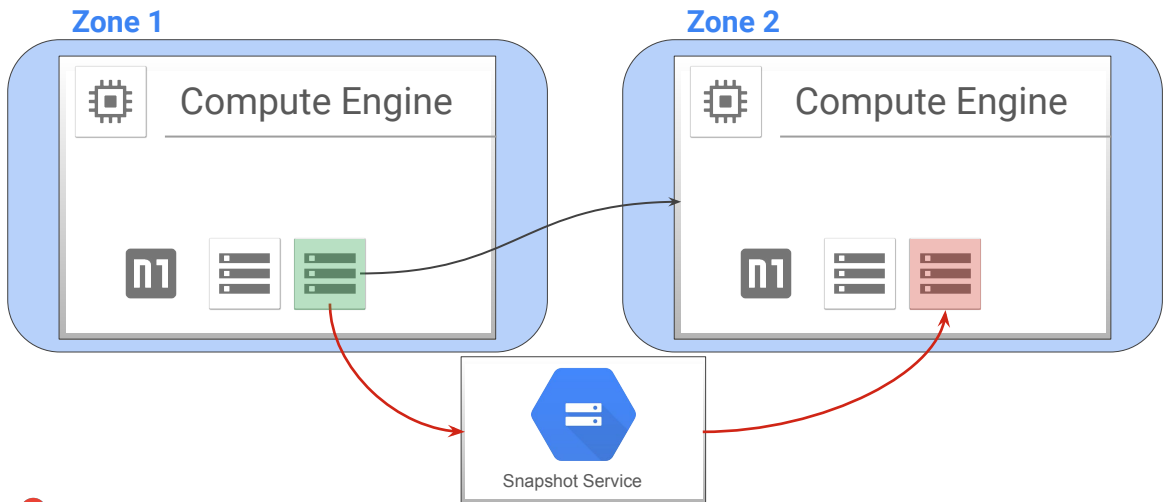
Speaking of snapshots, let's take a closer look at these.

Snapshot: Back up critical data



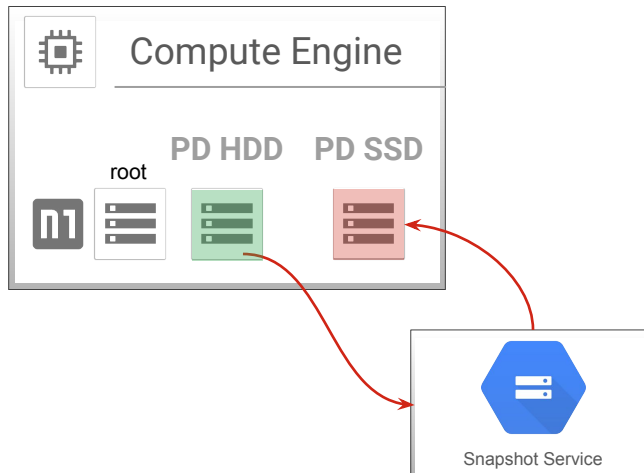
Snapshots have many use cases. For example, they can be used to back up critical data into a durable storage solution to meet application, availability, and recovery requirements. These snapshots are stored in Cloud Storage, which is covered later.

Snapshot: Migrate data between zones



Snapshots can also be used to migrate data between zones. I just discussed this when going over the manual process of moving an instance between two regions, but this can also be used to simply transfer data from one zone to another. For example, you might want to minimize latency by migrating data to a drive that can be locally attached in the zone where it is used.

Snapshot: Transfer to SSD to improve performance



Which brings me to another snapshot use case of transferring data to a different disk type. For example, if you want to improve disk performance, you could use a snapshot to transfer data from a standard HDD persistent disk to a SSD persistent disk.

Persistent disk snapshots

- Snapshot is not available for local SSD.
- Creates an *incremental* backup to Cloud Storage.
 - Not visible in *your* buckets; managed by the snapshot service.
 - Consider cron jobs for periodic incremental backup.
- Snapshots can be restored to a new persistent disk.
 - New disk can be in another region or zone in the same project.
 - Basis of VM migration: "moving" a VM to a new zone.
 - Snapshot doesn't back up VM metadata, tags, etc.



Now that I've covered some of the snapshot use cases, let's explore the concept of a disk snapshot.

First of all, this slide is titled persistent disk snapshots because snapshots are available only to persistent disks and not to local SSDs.

Snapshots are different from public images and custom images, which are used primarily to create instances or configure instance templates, in that snapshots are useful for periodic backup of the data on your persistent disks.

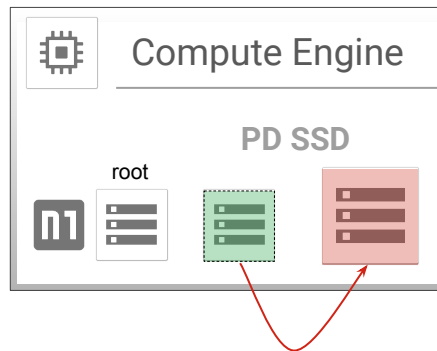
Snapshots are incremental and automatically compressed, so you can create regular snapshots on a persistent disk faster and at a much lower cost than if you regularly created a full image of the disk.

As we saw with the previous examples, snapshots can be restored to a new persistent disk, allowing for a move to a new zone.

To create a persistent disk snapshot, see the links section of this video:

[\[https://cloud.google.com/compute/docs/disks/create-snapshots\]](https://cloud.google.com/compute/docs/disks/create-snapshots)

Resize persistent disk



You can grow disks, but never shrink them!



Another common Compute Engine action is to resize your persistent disk. The added benefit of increasing storage capacity is to improve I/O performance. This can be achieved while the disk is attached to a running VM without having to create a snapshot.

Now, while you can grow disks in size, you can never shrink them, so keep this in mind.

Lab

Working with Virtual Machines



Let's get started with the second lab of this module.

In this lab, you'll be setting up an application server. Now this example happens to be a gaming application, but it applies to many other use cases. You will configure the VM and also add capacity for a production gaming system, and you will build the infrastructure that you need for production activities. This includes backups and graceful shutdown and restart services.

Lab review

Working with Virtual Machines



In this lab, you created a customized virtual machine instance by installing base software, which was a headless Java runtime environment and application software, specifically, a Minecraft game server.

You customized the VM by preparing and attaching a high-speed SSD, and you reserved a static external IP address so that the address would remain consistent.

Using that IP address, you then verified the availability of the gaming server online.

Next, you set up a backup system to back up the server's data to a Cloud Storage bucket, and you then tested that backup system. You then automated backups using cron.

Finally, you set up maintenance scripts using metadata for graceful startup and shutdown of the server.

Many of these techniques, including the script automation, can be adapted to administration of production servers in any application.

You can stay for a lab walkthrough, but remember that GCP's user interface can change, so your environment might look slightly different.

Review

Virtual Machines



In this module, we covered the different compute, image, and disk options within Compute Engine, along with some common actions. The two labs provided you with real world applications of most of the topics covered in this course.

Remember that there are many compute options to choose from. If a predefined machine type does not meet your needs, you can customize your own VM and you can even create a sole-tenant node. You can also install different public and custom images on the boot disks of your instances and you can attach more disks if needed.

Review

Essential Cloud Infrastructure: Foundation



Thank you for taking the Essential Cloud Infrastructure: Foundation course. I hope you have a better understanding of how to architect with Compute Engine, and I also hope that the demos and labs made you feel more comfortable with using the different GCP services that we covered.

Essential Cloud Infrastructure: Core Services

1. Cloud IAM
2. Data Storage Services
3. Resource Management
4. Resource Monitoring



Next, I recommend enrolling in the “Essential Cloud Infrastructure: Core Services” course of the “Architecting with Google Compute Engine” series.

In that course, we start by talking about Cloud IAM, and you will administer Identity and Access Management for resources.

Next, we’ll cover the different data storage services in GCP, and you will implement some of those services.

Then, we’ll go over resource management, where you will manage and examine billing of GCP resources.

Lastly, we’ll talk about resource monitoring, and you will monitor GCP resources using Stackdriver services.

Enjoy that course!