

# Compute Engine

Google Compute Engine lets you create and run virtual machines on Google infrastructure.

Google Compute Engine Concepts

## Features

- Large-scale IaaS workloads.
- High CPU, high memory, standard, and shared-core machine types.
- Persistent disks on standard (HDD), SSD, and local SSD.
- Resize disks and migrate instances with no downtime.
- Startup scripts and metadata.
- Robust networking:
  - Default and custom networks.
  - Firewall rules.
  - Regional HTTP(s) load balancing.
  - Network load balancing.
  - Subnetworks.
- Advanced APIs for auto-scaling and instance group management.
- Per-minute billing with sustained use discounts.
- Preemptible instances.
- High throughput to storage at no extra cost.
- Custom machine types allow you to only pay for what you need.
- Import VM uses Cloud Endure. Super easy.

## Points of Interest

- Windows virtual machines require you to set a username and password.
- Configuration management options include:
  - Startup script for Compute instances.
  - Google Cloud Deployment Manager.
  - Open-source tools such as Puppet, Chef, Salt, and Ansible.

## Shutdown Script

- Resetting an instance does not run the shutdown script.
- Restart, reboot, stop, shutdown, delete will allow approximately 90 seconds for the shutdown script to run.
- Preemptible machines will allow only 30 seconds for the shutdown script to run.

## Changing Machine Type

- Set the boot disk to not be deleted when the virtual machine is deleted.
- Delete the virtual machine instance.
- Create a new instance with the new machine type settings.
- Select the boot disk from the existing disks.

## Preemptible Machines

- Lower price for interruptible service (up to 80%),
- May be terminated at any time:
  - No charge if within 10 min.

- 24 hours max run time.
- 30 second termination notification.
- No live migrate.
- No auto restart.
- Can request a CPU quota be split between regular and preemption.

## Disks

- HDD or SSD.
- Live disk resizing (bigger, not smaller).
- Supports attachment to multiple VMs in read only mode.
- Automatic Checksums.
- Automatic Encryption (can supply own keys).
- Persistent Disk:
  - Network storage appearing as a block device.
  - Attached to the VM through the network interface.
  - Durable storage.
  - Bound to zone.
  - Bootable.
  - Snapshots.
- Local SSD Disk:
  - Physically attached to VM.
  - Not available on shared core.
  - Faster than Persistent disk.
  - Ephemeral: data survives a restart but not a stop or terminate.
  - 3TB (375GB \* 8).
- RAM Disk:
  - tmpfs
  - Faster than Local SSD, slower than memory.
  - Very volatile.
  - RAM required so larger machine type needed.

	Persistent HDD	Persistent SSD	Local SSD	RAM
Redundancy	Yes	Yes	No	No
Encryption	Yes	Yes	Yes	N/A
Snapshotting	Yes	Yes	No	No
Bootable	Yes	Yes	No	No
Use Case	Bulk File	Random IOPS	High IOPS Low Lat	Low Lat

No of Cores	Disk Limit
Shared Core	16
1 Core	32
2-4 Cores	64
8 or more	128

## Moving a VM to a new zone

- Manual process:
  - Snapshot all persistent disks on the source VM.
  - Create new persistent disks in destination zone restored from snapshots.
  - Promote ephemeral external IP to static external IP.
  - Create new VM in the destination zone and attach new persistent disks.

- Assign static IP to new VM, demote to ephemeral.
  - Update references to VM.
  - Delete the snapshots, original disks, and original VM.
- Automated process:
  - `gcloud compute instances move`
  - Update references to VM.

## Performance Management

- Region choice affects machine type and CPU architecture options.
- 1 vCPU is equal to 1 hyperthreaded core.
- 2 vCPUs is equal to 1 physical core.
- Network throughput is 2Gbps per vCPU up to 16Gbps for 8 vCPUs.
- Disk throughput is tied to the network throughput.
- Disk IOPs is tied to the disk size.
- Local SSD is ephemeral and ties the instance to the hardware.

## Availability Policy

- Preemptibility: On or Off.
- Automatic restart: On or Off.
- On host maintenance: Migrate or Terminate.
- GPUs will prevent migration on host maintenance.
- Local SSD disks will prevent migration on host maintenance.

## Pricing

- Per minute billing with a 10 minute minimum.
- Sustained use discounts.
- Preemptible instances:
  - Live at most 24 hours.
  - Can be pre-empted with a 30 second notification via API.
  - Up to 80% discount.
- Custom machine types will have custom pricing.
- Recommendation engine notifies of under utilized instances.
- Committed use discounts (1 year or 3 years).
- Inferred instance discounts: Usage of VMs of the same machine type in the same zone are combined as if they were one machine.
- No charge for stopped instances apart from attached disks and IPs.
- Same charge for different CPU architectures. Choose your region wisely.

## Autoscaling

- Automatically scale the number of instances in the managed instance group based on workload.
- Can reduce costs by shutting down instances when not required.
- Create one autoscaler per managed instance group.
- Support both zone-based managed instance groups or regional managed instance groups.
- It is fast responding typically within a 1 minute moving window.
- Policies include Max and Min number of replicas.
- Policy options:
  - Average CPU utilization.
  - HTTP load balancing serving capacity (backend service definition for Max CPU or Max req/sec/instance.)
  - Stackdriver standard and custom metrics.

- Google Cloud Pub/Sub queuing workload.
- Supports up to 5 policies.
- Use Stackdriver custom metrics to autoscale more accurately to your application workload.

### **Autoscaler Configuration**

1. Create instance template (startup scripts, shutdown scripts, software, logging).
2. Create managed instance group.
3. create autoscaler.
4. Optionally, define multiple policies.