

# Compute Engine

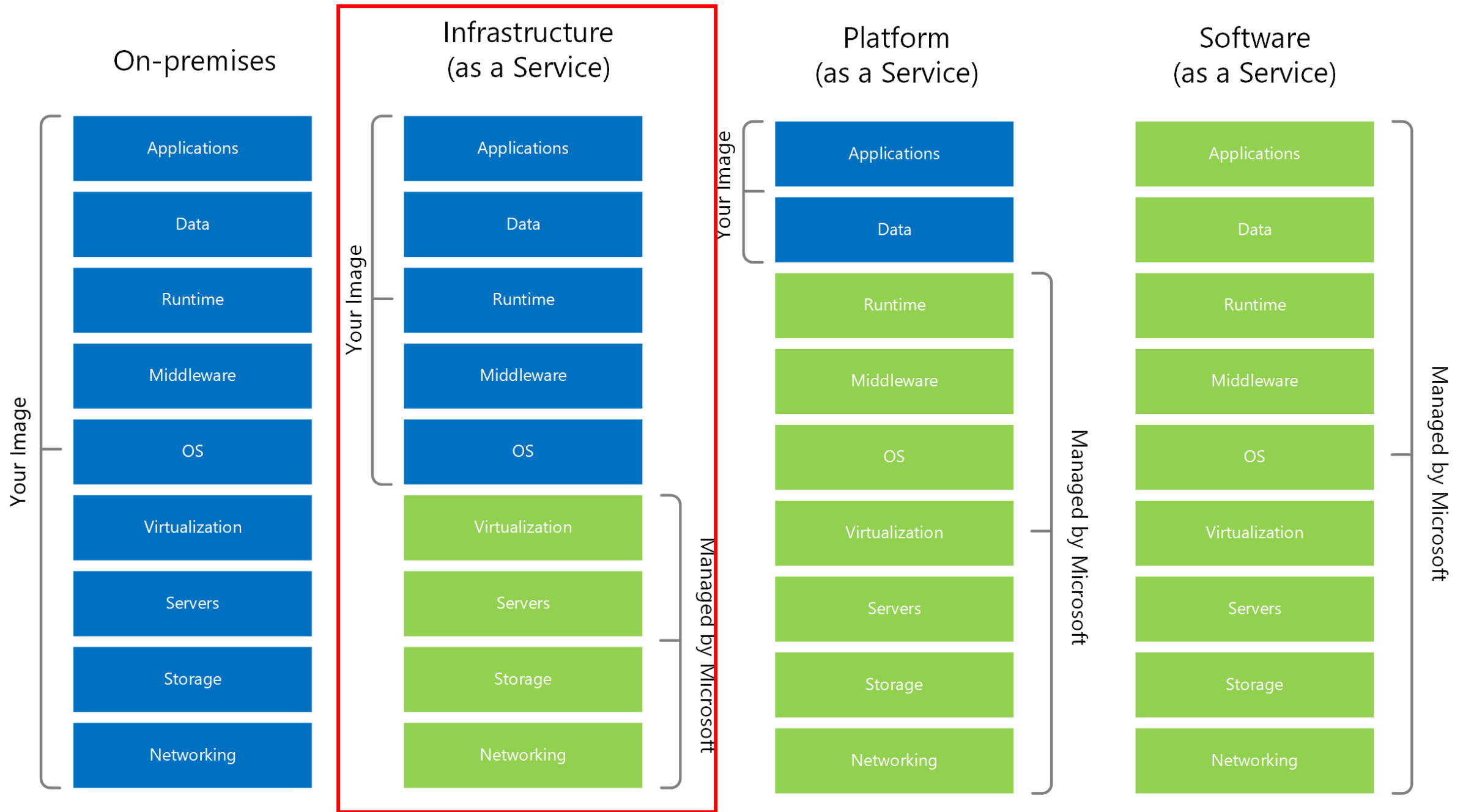
Memi Lavi  
[www.memilavi.com](http://www.memilavi.com)



# Compute Engine

---

- Google Cloud service for creating and managing virtual machines
- Offers wide variety of VM-related services
- IaaS



# Virtual Machines

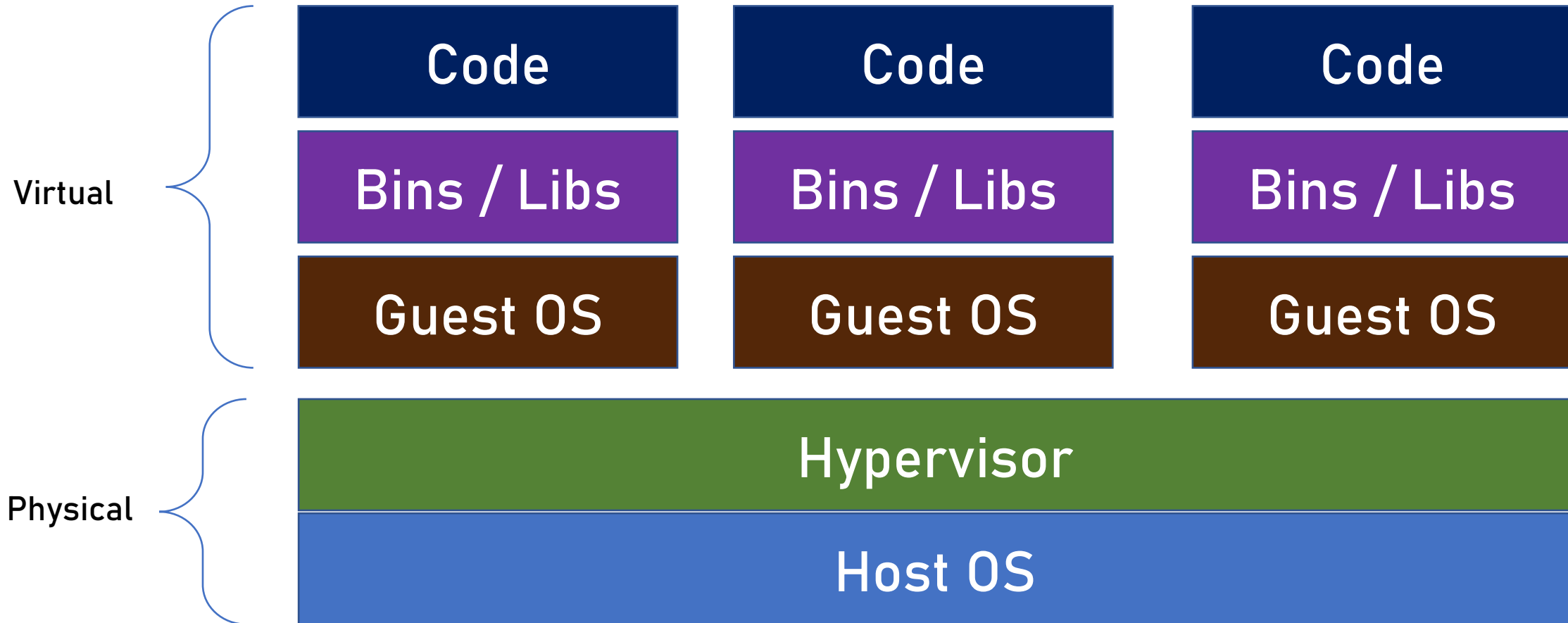
---

- A virtual (=not real) server running on a physical (=real) server
- Allows creating new servers extremely quick
- Based on existing resources of the physical server
- From the user's point of view – a regular server, nothing new
- SLA of a single machine: 99.5%

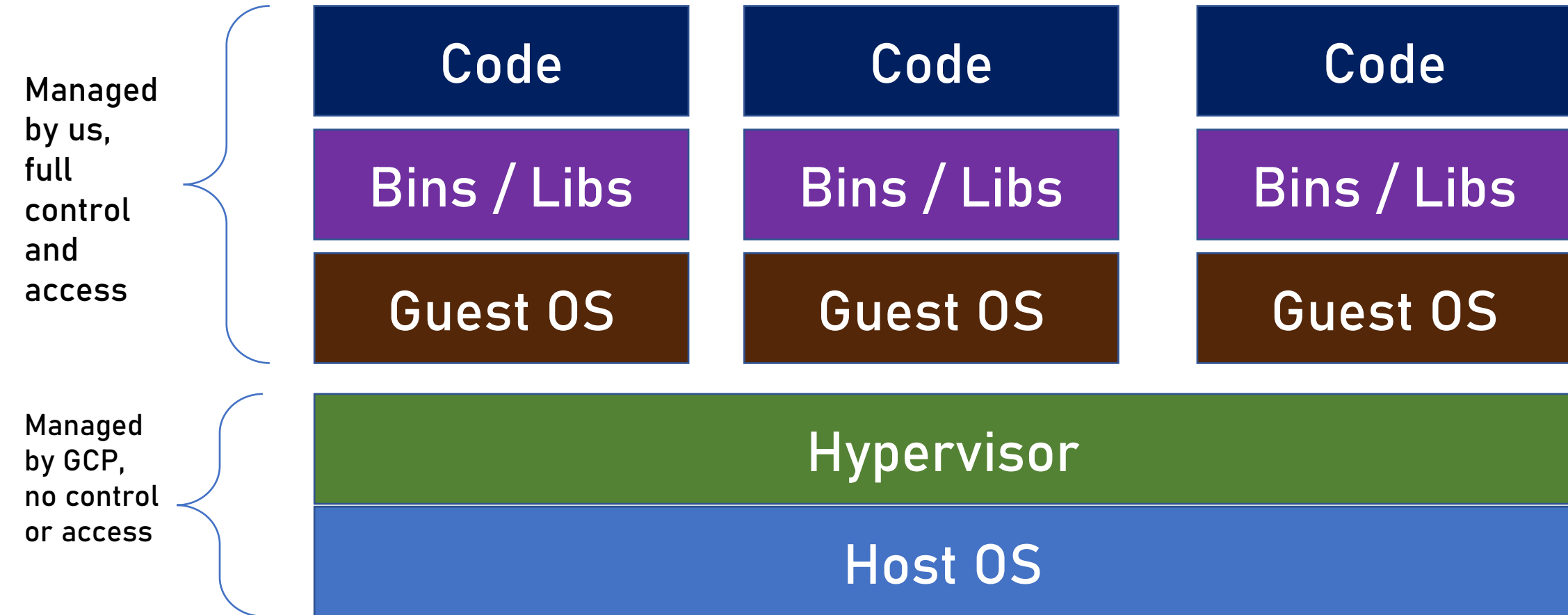
# Virtual Machines Architecture

VM Density = No. of VMs per Host

In this case - 3



# Virtual Machines In GCP



# Virtual Machines in GCP

- Steps for creating VM in GCP:

- Select the location

- Select the image (OS + Pre-Installed software)

**Don't forget to check the price!**

- Select the size

- Customize as needed

- That's it, basically....

# VM Families

---

- VMs in GCP are divided to families
- Each family targets specific type of workload
- Families are further classified by series and generation (usually newer hardware)
- Always make sure you select the appropriate VM family for your needs



# Workload Types

- VM families are classified to the following workload types:

## General Purpose

Web servers, microservices, virtual desktops, databases, etc.

## Compute Optimized

High-performance computing (HPC), game servers, media transcoding, etc.

## Memory Optimized

SAP HANA databases, in-memory data stores (Redis), etc.

## Accelerator Optimized

ML training, massively parallelized computation, deep learning etc.

# General Purpose

---

- VM series:

E2

Low traffic web server, back office apps, virtual desktops

N2, N2D, N1

Medium traffic web servers, microservices, BI

C3

High traffic web servers, databases, in-memory cache

TAU T2D, T2A

ARM architecture, Cost effective, scale out workloads

# General Purpose Hardware

Series	# of CPU	RAM	CPU Family
E2	Up to 32	Up to 128GB	Intel and AMD EPYC
N1	Up to 96	Up to 624GB	Intel, multiple families
N2	Up to 128	Up to 864GB	Intel Ice Lake and Cascade Lake
N2D	Up to 224	Up to 864GB	AMD EPYC Milan and AMD EPYC Rome
C3	Up to 176	Up to 1.4TB	Intel Sapphire Rapids
TAU T2A	Up to 48	Up to 192GB	Ampere Altra Arm
TAU T2D	Up to 60	Up to 240GB	AMD EPYC Milan

# Compute Optimized

---

- VM series:

H3

HPC workload, scientific and engineering computing

C2, C2D

Gaming, ad serving, media serving, AI/ML

# Compute Optimized Hardware

Series	# of CPU	RAM	CPU Family
C2	Up to 60	Up to 240GB	Intel Cascade Lake
C2D	Up to 112	Up to 896GB	AMD EPYC Milan
H3	88	352	Intel Sapphire Rapids

# Memory Optimized

---

- VM series:

M1, M2, M3

SAP HANA, MS SQL, memory intensive apps

# Memory Optimized Hardware

Series	# of CPU	RAM	CPU Family
M1	Up to 160	Up to 4TB	Intel Cascade Lake
M2	Up to 416	Up to 12TB	Intel Cascade Lake
M3	Up to 128	Up to 4TB	Intel Ice Lake

# Accelerator Optimized

---

- VM series:

A2, A3

HPC, ML Training

G2

Video transcoding, Remote visualization



# Memory Optimized Hardware

Series	# of CPU	RAM	GPU Family
A2	Up to 96	Up to 1.4TB	NVIDIA A100
A3	Up to 96	Up to 2TB	NVIDIA H100
G2	Up to 96	Up to 432GB	NVIDIA L4

# GPU

---

- Some VM series can use GPU in addition to CPU
- The A2 and G2 series have built in support for GPU and no need to add it later
- GPU can also be added to N1 machines
- All GPUs are NVIDIA based

# Custom Machines

---

- If none of the predefined machine types fit your need you can create your own custom machine
- Custom machines can be created in the E2, N2, N2D and N1 series
- You define the number of CPU and the size of RAM
- For example: 32CPU, 16GB RAM
- Not all combinations are valid

# Reducing Cost of VM

---

- VM Instances can quickly become quite expensive
- There are various techniques to reduce the cost of a VM
- Can be combined
- Choose wisely

# Reducing Cost of VM

---

Scheduled Shutdown

Spot Instances

Committed Use Discounts

Disk Types

# Scheduled Shutdown

---

- VM operations can be scheduled
- Useful for VMs that are not needed 24X7
  - Test / Dev
  - Batch processing
- Instances can be shutdown when not in use, saving costs
- Done using Instance Schedules

# Spot Instances

---

- Great way to save costs
- Up to 91% discount
- Instances that can be stopped and removed by Google at any time to use their resources
- Use for batch processing, data analytics etc.
- Don't use for web apps and systems that need to be accessible

# Committed Use Discounts

---

- Get deep discount for committing to use resources
  - Up to 70% discount
- You pay the committed price regardless of the actual use
- Commitment is for 1 or 3 years
- Can commit to specific hardware or license



# Disk Types

---

- When creating VM instance it's important to select the best disk type for it
- Disk type affect the reliability, speed and cost

# Disk Type Families

- Two major disk types:

## Local SSD

- Attached to the physical host of the VM
- Provides the best performance
- Ephemeral – when the VM shuts down data is lost
- Best for temporary data
- Size is always 375GB
- Can attach multiple disks per VM
- Price depends on region, usually 0.08\$ / GB
  - ~30\$ / month

## Persistent

- Durable network storage
- Zonal or Regional
- Can be detached and moved from a VM
- Auto scale

# Persistent Disk Types

## Balanced

- SSD
- Balances performance and cost
- Best for general purpose systems

## Performance

- SSD
- For high-performance databases
- Designed for single-digit millisecond latencies

## Standard

- HDD
- Large data processing

## Extreme

- SSD
- Consistent high performance
- Provision your target IOPS

# VM Metadata

---

- Every VM stores metadata in a special metadata server
- Metadata can be accessed using specialized HTTP requests from the VM
- Contains various data about the VM and the project
- Can be customized per VM
- Can be set by apps running on the machine

# VM Metadata

---

- Useful for:
  - Get data about the VM, ie. external IP
  - Use data in startup script
  - Learn about maintenance schedule
  - And more...

# Accessing VM Metadata

- From within the VM:

- Access root URL:

<http://metadata.google.internal/computeMetadata/v1>

- Specify header:

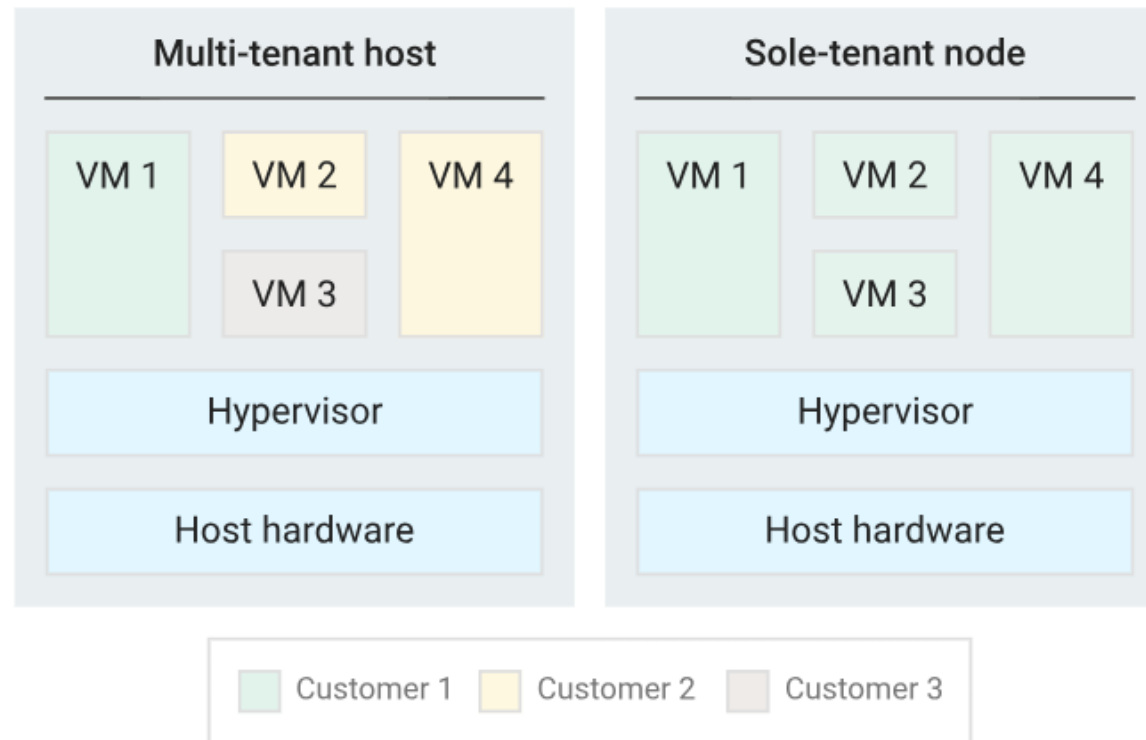
Metadata-Flavor: Google

# Sole Tenancy

---

- By default, VM instances are placed on physical servers that host VM instances of other GCP customers
- You can have a dedicated physical server for your VMs
- This is called “Sole Tenant Nodes”
- Quite expensive

# Sole Tenancy





# Instance Templates

---

- Organizations usually need many VMs
- These VMs should have identical configuration
  - Family, disk, image, etc.
- Creating instances manually is error prone and slow
- Instance Templates define a VM configuration that can be used to create instances quickly and easily

# Instance Templates

- Instance Templates are:

Global

But can include regional / zonal configuration, for example – disk in a specific region

Free

You can create as many templates as you want. You'll pay only for a VM instance you'll create from them

Immutable

You cannot change Instance Template once it was created, to ensure consistency of the VM Instances created from it

# Snapshots

---

- Exact copy of a disk in a specific point in time
- Used for backup
- Incremental
  - Second snapshot contains the difference from the first one
  - Saves space and cost
- Can be scheduled

# Snapshot Types

---

## Standard

Geo-redundant, great for backup against local, zonal, regional outage. Can be scheduled.

## Archive

Geo-redundant. Rarely accessed, cannot be scheduled, cost effective.

## Instant

Local, for use when need a quick restore to a new disk.

# Images

---

- Store configuration and data of a machine
- Can be used for backup, maintenance and cloning
- Can be used to backup multiple disks at a time
  - A snapshot backs up only a single disk
- Stored in a cloud storage bucket

# Images

---

- Machine image stores:
  - Description
  - Machine type
  - Instance metadata
  - Labels
  - Network tags
  - Policies
  - Disks data

# Images

---

- Machine image does not include:
  - Local SSD data
  - Memory
  - Name and IP address

# Image Types

---

Machine

Contains all the data of the machine, including its disks data

Disk

Copy of a disk, can be used as a boot disk for a new machine



# Image Cost

---

- Based on the image size and its storage location
- We'll learn about storage costs later

# Instance Groups

---

- A collection of VM instances that's managed as a single entity
- Great for systems that require more than a single VM
- If deployed in multiple zones, SLA is 99.99%

# Instance Group Types

Managed

Use for:

- Web apps
- Databases
- Batch processing

Unmanaged

Use for:

- Heterogenous workloads

Always try  
to use this

The diagram consists of two purple rectangular boxes at the top, labeled 'Managed' on the left and 'Unmanaged' on the right. Below the 'Managed' box is a list of use cases: 'Web apps', 'Databases', and 'Batch processing'. Below the 'Unmanaged' box is a single use case: 'Heterogenous workloads'. At the bottom center, there is a hand-drawn oval containing the text 'Always try to use this'. A brown arrow originates from the bottom of this oval and points diagonally upwards to the bottom edge of the 'Managed' box.

# Managed Instance Group (MIG)

---

## High availability

- Auto-healing instances by recreating them
- Support for app-based health checks
- Regional
- Load balancing support

## Scalability

- Auto scale
- Add and remove instances as needed

## Automatic Updates

- Support for various deployment scenarios
- Full control on speed and scope of deployment

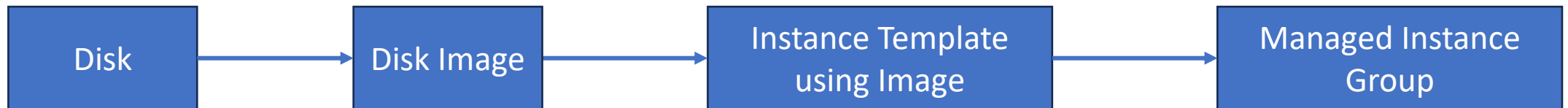
# Unmanaged Instance Groups

---

- Use when you need full control on the instances
- Support for heterogenous instances

# Creating Instance Groups

- Instance Groups are created from Instance Templates
- Instance Templates can be based on Images



# VM Manager

---

- When having many VMs in a project it's not easy to manage them
- We need to:
  - Patch the VMs
  - Manage inventory
  - Update software

# VM Manager

---

- VM manager is a suite of tools in GCP that help manage VM fleet
- Three main services:

OS Patch Management

Apply on-demand and scheduled OS patches

OS Inventory Management

Collect and review operating system information

OS Configuration Management

Install, remove and auto-update software packages



# VM Manager

---

- Should be enabled
  - Per project
  - - Or -
  - Per VM
- Free up to 100 VMs

# GCP Architecture Diagram

---

- When designing architecture for GCP apps it's a good idea to use GCP symbols in the diagram
- There are hundreds of them...
- Follow GCP best practices for that

# Download GCP Icons and Slides

---

- <https://cloud.google.com/icons>

## Architecture: ReadIt Cloud System

