

Cloud Computing Cloud Applications

CIS437

Erik Fredericks // frederer@gvsu.edu

A question

What makes a cloud application different from a normal application?

And further

Are there different types of cloud applications?

Limitations of "local" apps

...remember the advantages of cloud computing?

- Fault tolerance
- Scalability
- Networking
- Federation
- Computing at scale
- ...



First, some things to consider before we dip into *aaS

APIs

Availability zones

Fault tolerance

Migration

Monitoring

Federation

Elasticity

Architectures

i.e., important things to consider in your cloud app

APIs

What is an API (other than something we'll talk about more when we go over serverless/microservices)?

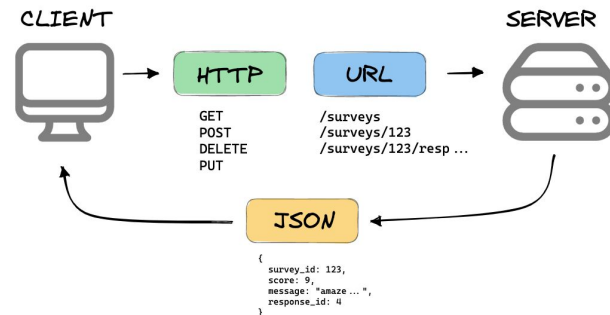
A method for accessing services programmatically

- Could be a human or a program!

Why important to cloud?

- Ability to write code to automate *everything* in the cloud!

WHAT IS A REST API?



mannhowie.com

```
$ aws ec2 run-instances \
  --image-id ami-1a2b3c4d \
  --count 1 \
  --instance-type c3.large \
  --key-name MyKeyPair \
  --security-groups MySecurityGroup
```

(a)

```
$ openstack server create --flavor 1 --image 397e713c-b95b-4186-ad46-6126863ea0a9 \
  --security-group default --key-name KeyPair01 --user-data cloudinit.file \
  myCirrosServer
```

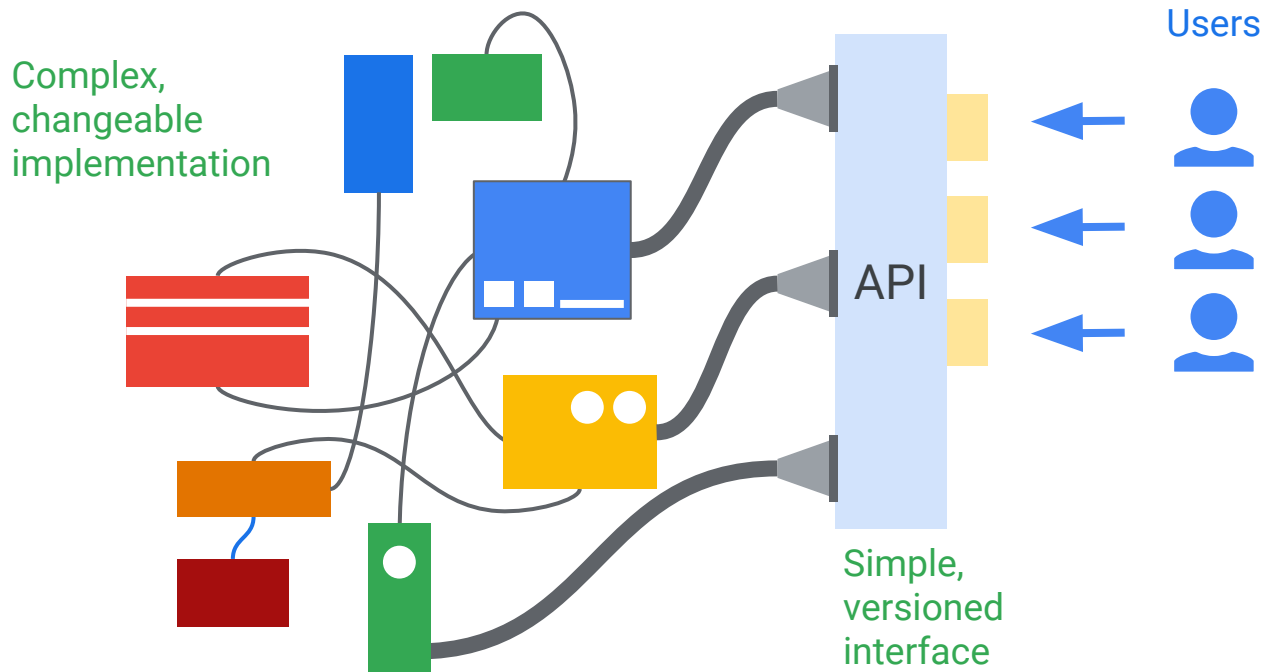
(b)

```
$ gcloud compute instances create "my-new-instance" \
  --zone="us-west1-b" \
  --image-family="tf-latest-cu92" \
  --image-project="deeplearning-platform-release" \
  --maintenance-policy=TERMINATE \
  --accelerator="type=nvidia-tesla-v100,count=8" \
  --machine-type="n1-standard-8" \
  --boot-disk-size=120GB \
  --metadata="install-nvidia-driver=True"
```

(c)

Figure 3.1: How to run a new VM by using the command-line client software provided by (a) AWS, (b) OpenStack and (c) GCP respectively.

APIs hide the details and enforce contracts



Availability zones

or, where does your provider have service?

Compute/serve near your client

- Faster delivery / quality of service

Regions can have multiple zones as well!



33 launched Regions
each with multiple Availability Zones

105 Availability Zones

600+ CloudFront POPs
and 13 Regional edge caches

AWS Global Infrastructure Map

The AWS Cloud spans 105 Availability Zones within 33 geographic regions, with announced plans for 21 more Availability Zones and seven more AWS Regions in Malaysia, Mexico, New Zealand, the Kingdom of Saudi Arabia, Thailand, Taiwan, and the AWS European Sovereign Cloud.



List view

41 Local Zones
29 Wavelength Zones
for ultralow-latency applications

245 countries and
territories served

135 Direct Connect
locations

<https://cloud.google.com/about/locations#lightbox-regions-map>
<https://aws.amazon.com/about-aws/global-infrastructure/?p=ngi&loc=1>

Fault tolerance

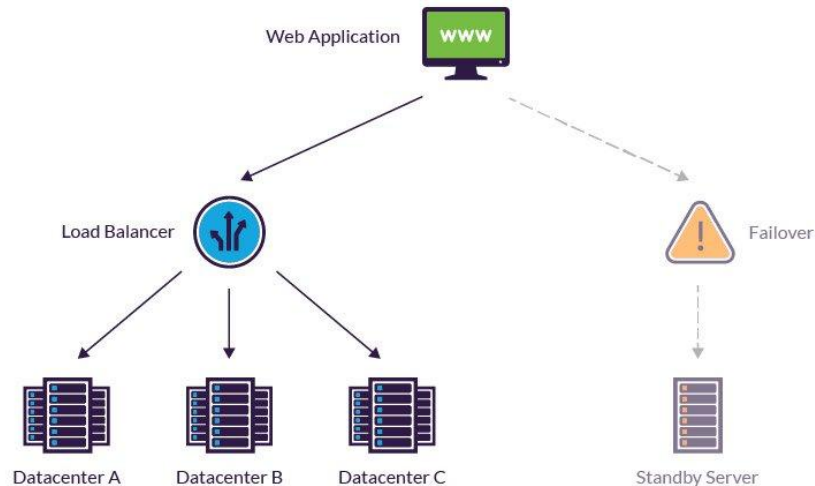
There are a *lot* of resources available

- What do we do if one fails?
 - e.g., a virtual machine crashes, or a zone goes down

Services need to be configured to *fail nicely*

- E.g., a mirror picks up the slack (failover)
 - Another instance is automatically spun up
 - etc.

Could be provided automatically by the provider or configured by the designer



Migration

May need to move from one service to another


- E.g., one VM needs to be moved to another
 - Because of availability, new services, etc.

Hot (live) migration

- Services still "on" while migration happens

Cold migration

- Services "off" while migration happens



Why one over the other?

Monitoring

KEEPIN' AN EYE ON THINGS

- Why on earth would we want to monitor our services?
- What can we do if things *go wrong*?

Google Developer

< Projects

Cloud Logs Test Pr...

Overview

Permissions

Billing & settings

APIs & auth

APIs

Credentials

Consent screen

Push

Monitoring

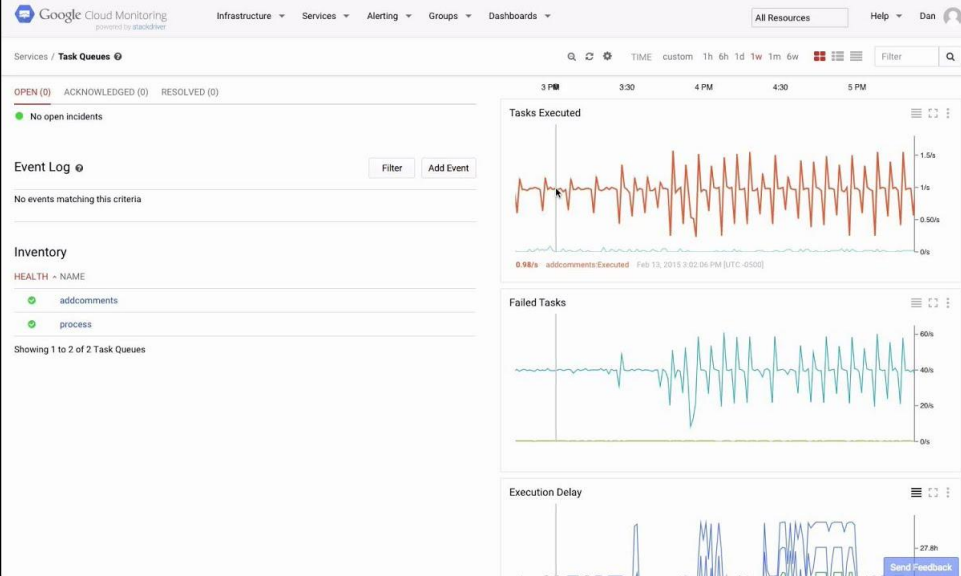
Traces

Logs

Dashboards & alerts

Source Code

Compute



Filter by label or text search

Compute Engine

All resource types

All resource IDs

All logs

Any log level

2015-03-16

▶	✖	mju-ubuntu	16:17:01.000	Mar 16 20:17:01	mju-ubuntu	CRON[12485]: (root) CMD (c
▶	✖	logging-docs	16:17:01.000	Mar 16 20:17:01	logging-docs	/USR/SBIN/CRON[15611]: (z
▶	✖	logging-docs	17:17:01.000	Mar 16 21:17:01	logging-docs	/USR/SBIN/CRON[16577]: (z
▶	✖	mju-ubuntu	18:17:01.000	Mar 16 22:17:01	mju-ubuntu	CRON[14487]: (root) CMD (c
▶	✖	logging-docs	18:17:01.000	Mar 16 22:17:01	logging-docs	/USR/SBIN/CRON[17530]: (z
▶	✖	mju-ubuntu	19:17:01.000	Mar 16 23:17:01	mju-ubuntu	CRON[15595]: (root) CMD (c
▶	✖	logging-docs	19:17:01.000	Mar 16 23:17:01	logging-docs	/USR/SBIN/CRON[18660]: (z
▶	✖	mju-ubuntu	20:17:01.000	Mar 17 00:17:01	mju-ubuntu	CRON[16598]: (root) CMD (c
▶	✖	logging-docs	20:17:01.000	Mar 17 00:17:01	logging-docs	/USR/SBIN/CRON[19754]: (z
▶	✖	mju-ubuntu	21:17:01.000	Mar 17 01:17:01	mju-ubuntu	CRON[17561]: (root) CMD (c
▶	✖	logging-docs	21:17:01.000	Mar 17 01:17:01	logging-docs	/USR/SBIN/CRON[20717]: (z

Federation

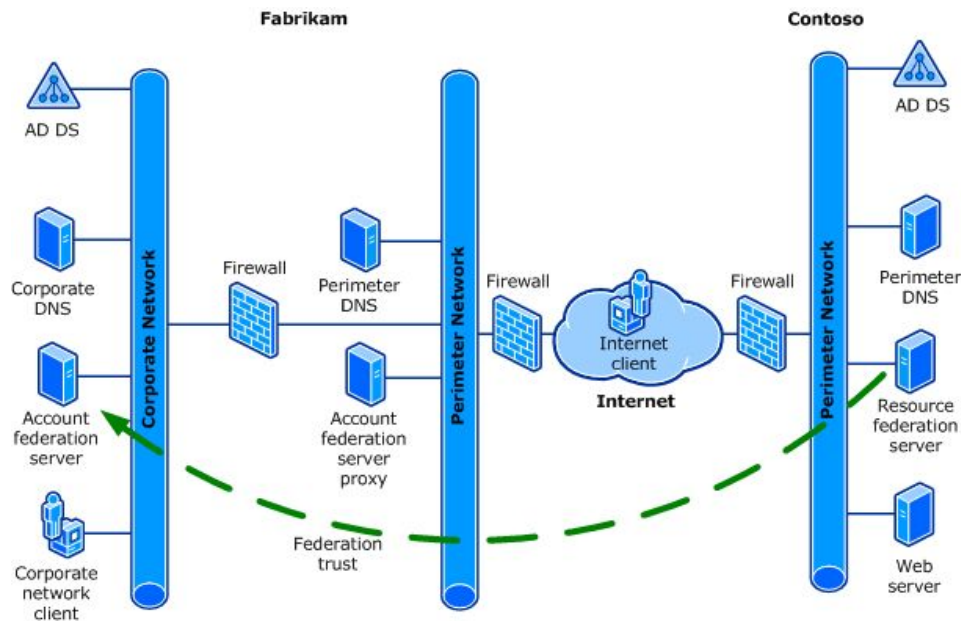
Interoperability of different providers

Think:

- Microsoft Active Directory
- Mastodon / Lemmy

Is this ... something that most cloud providers allow?

What do you think "lock-in" means?



Elasticity

This is what gives us scalability

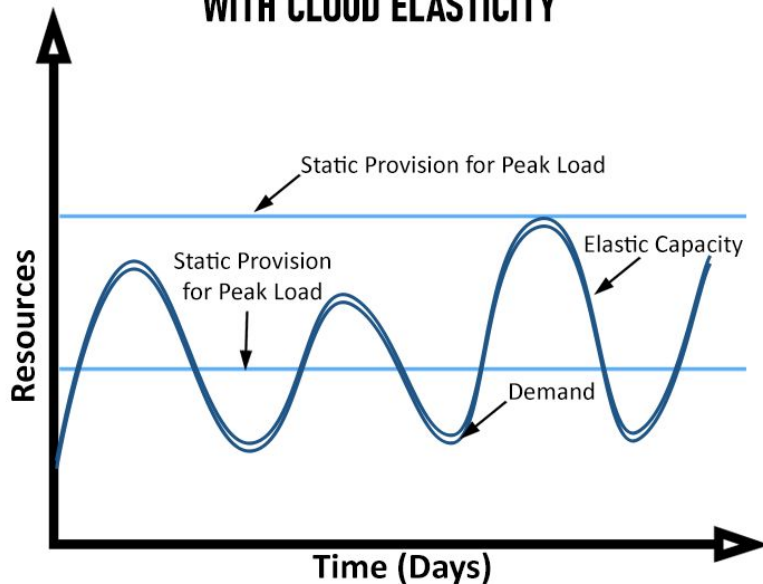
- Detect if workloads are too high and scale up/down as needed
- Relies on monitoring and automated triggering!

Think of your programs you've made so far though...

- How in the world could you automatically scale them up/down *at run time*?



COMPARISON OF STATIC CAPACITY WITH CLOUD ELASTICITY



Architecture

"Common"
(from the book)

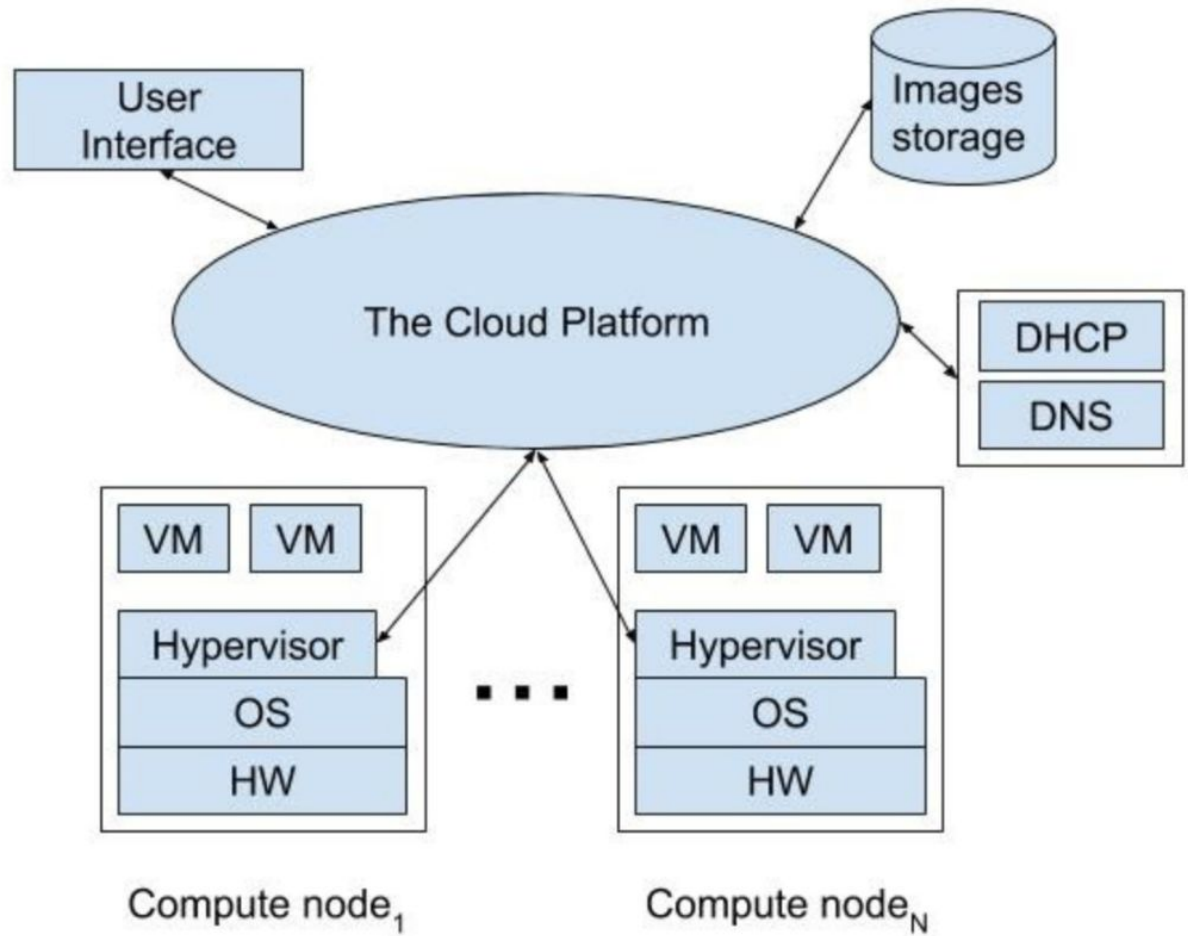


Figure 3.2: The common cloud platform architecture.

Many different types, but typically **three tiers** of apps

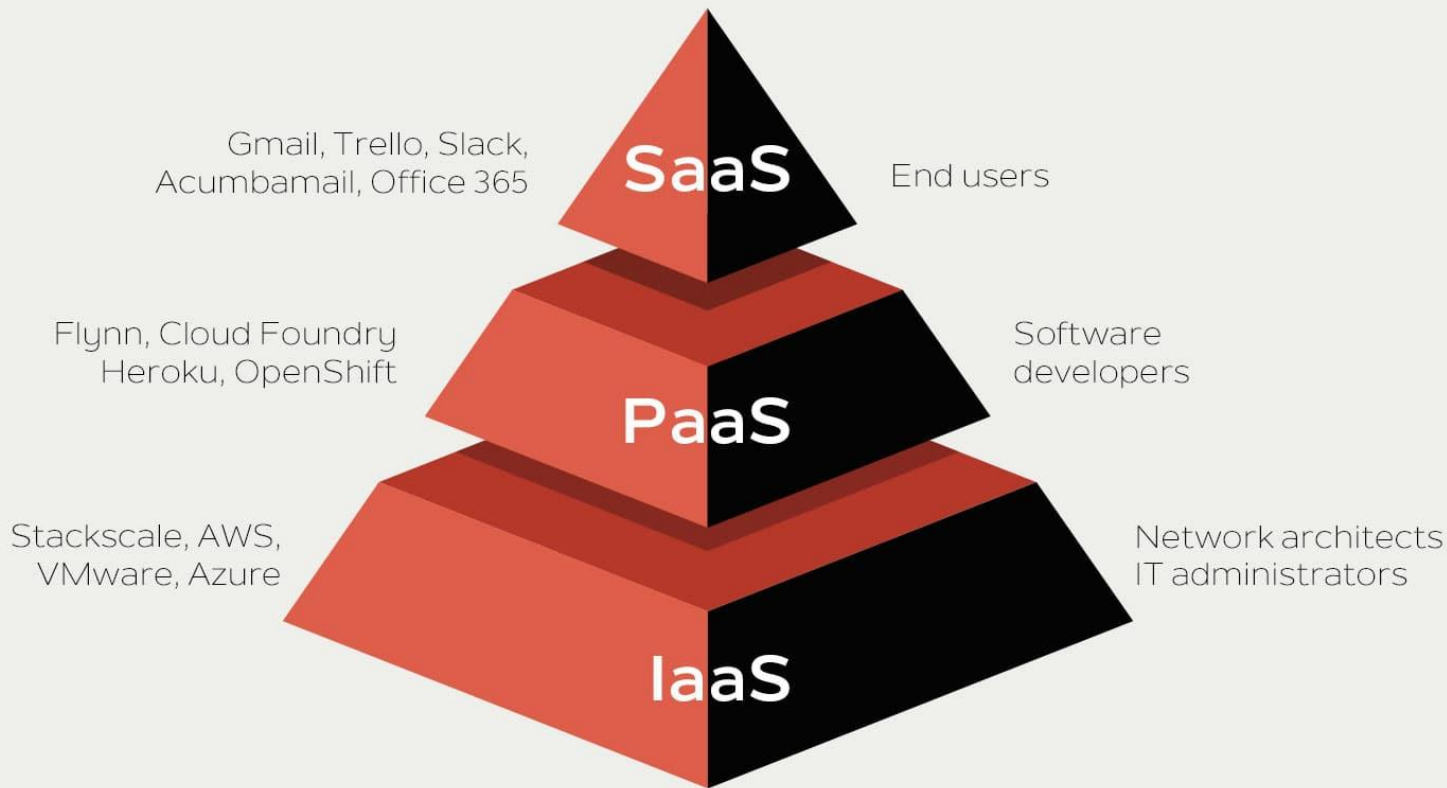
* as a service

- Infrastructure as a service (IaaS)
- Platform as a service (PaaS)
- Software as a service (SaaS)

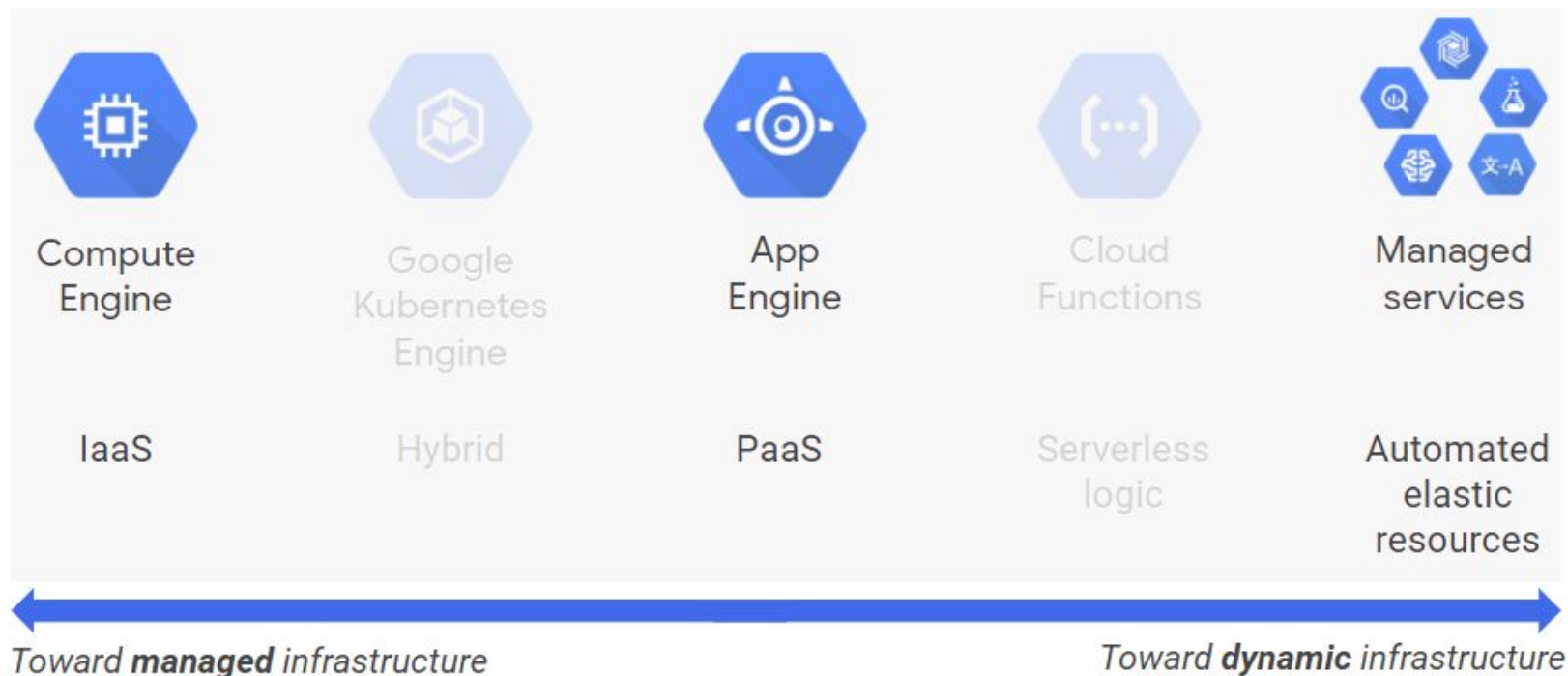
Interestingly, there will typically be some cloud product that aligns with each of these

- Or some combination

Cloud service models



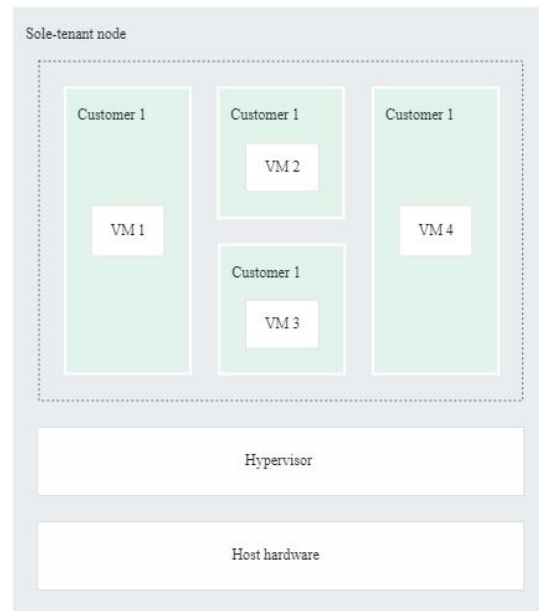
Google Cloud provides a variety of choices



IaaS

Typically **virtual machines**

- Meaning, the infrastructure is virtualized
 - Don't need a server "on-prem"



IaaS

Things to consider for the VM:

- Machine type/specs
- Operating system
- Who has access
- Fresh install vs. template

Other considerations:

- Firewalls (internal and external!)
- Scaling (elasticity)
- Location
- Price!
 - Use vs. storage



What are some sample use cases?

Autoscaling Google Compute Engine

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

Uses "managed instance groups"

- Auto add/delete instance from groups as needed

Requires:

- Autoscaling service agent
- A policy (metrics to look for)
 - CPU usage
 - Load balancing
 - Etc.

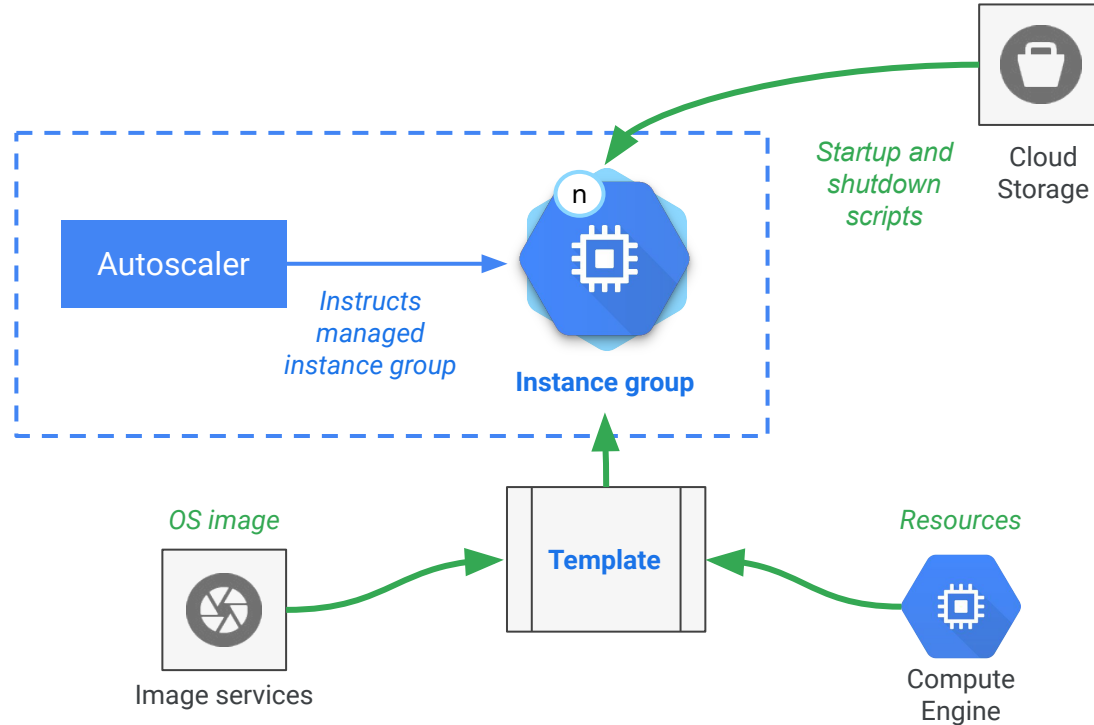
Create a managed instance group (MIG)

A [managed instance group \(MIG\)](#) is a group of virtual machine (VM) instances that you control as a single entity. MIGs support features such as autohealing, autoscaling, load balancing, multiple zone coverage, and stateful workloads.

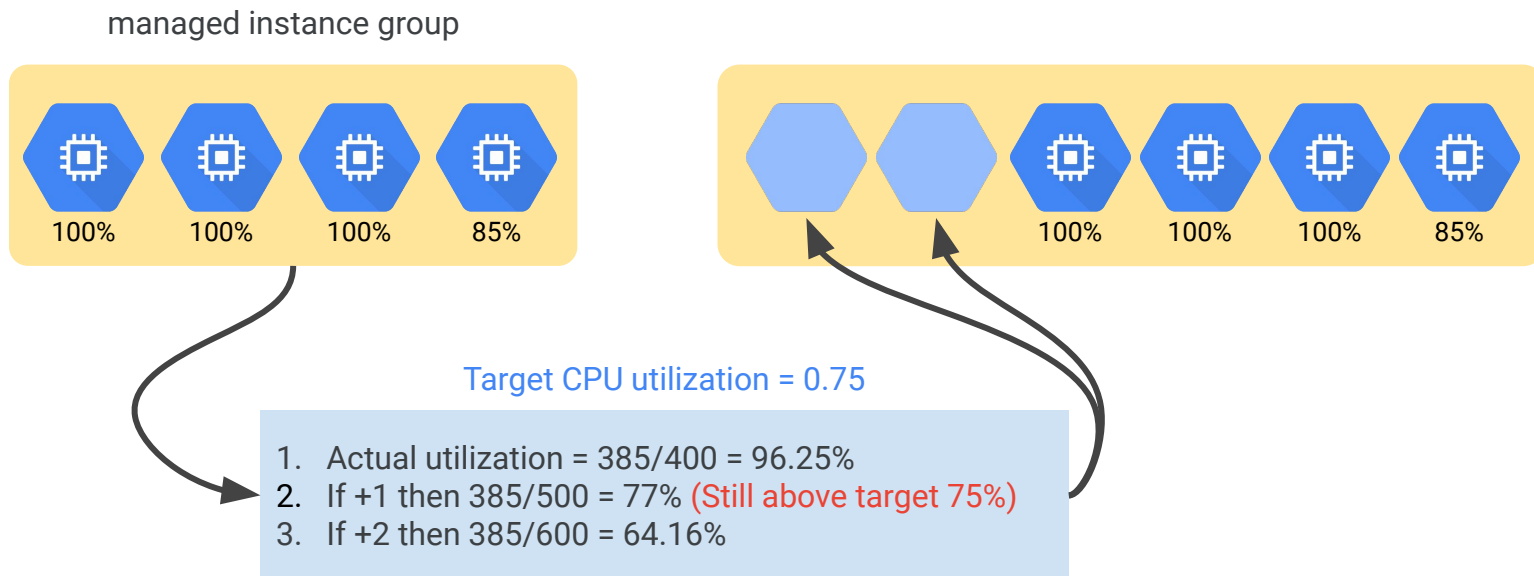


Quickstart
tutorial in VM
creation

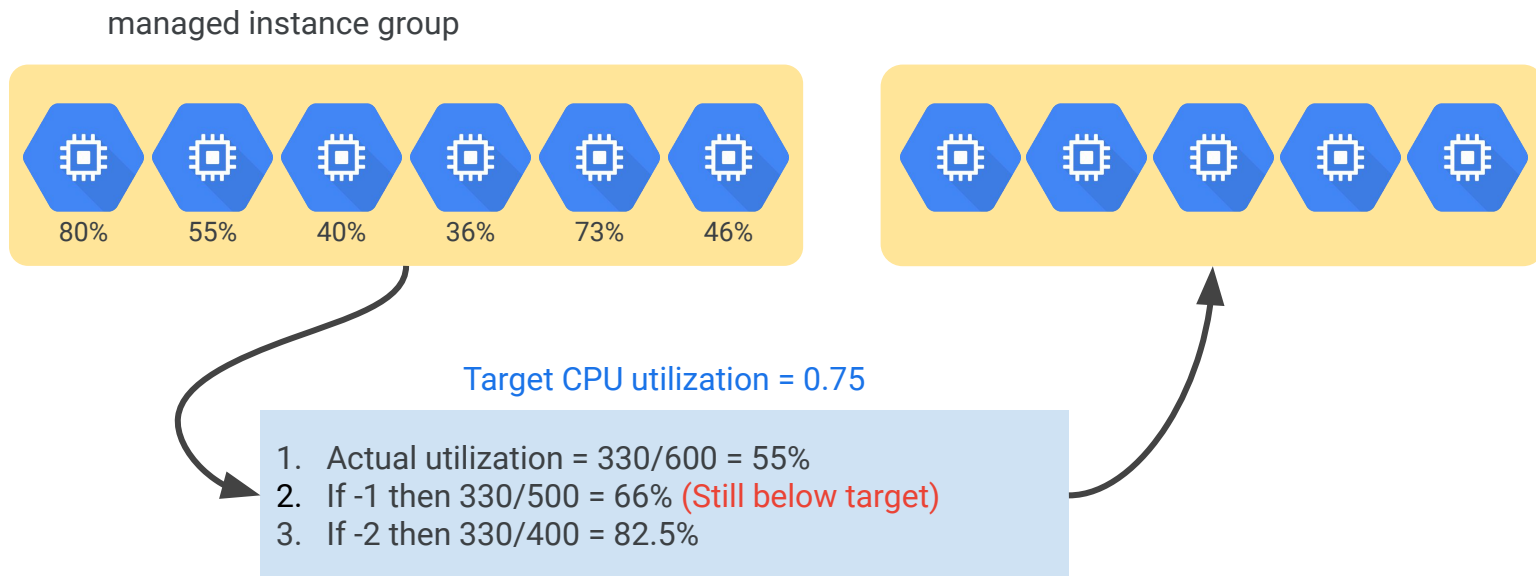
How autoscaling works



Scale-out policy decision



Scale-in policy decision



Autoscale demo

<https://codelabs.developers.google.com/codelabs/hpc-slurm-on-gcp#0>

~~(Interestingly, Slurm is what runs on GVSU's HPC environment)~~

~~— And many others...~~

~~Notes:~~

~~— Create new project~~

~~— Pull this instead: <https://github.com/GoogleCloudPlatform/slurm-gcp>~~

~~— Look in here for the tfvars:~~

~~slurm-gcp/terraform/slurm_cluster/examples/slurm_cluster/simple_cloud~~

aaaaaaAAAAAAAAAAAAAAAAAAAAHHHH

Well Slurm is at least fun let's try THAT out

<https://cloud.google.com/cluster-toolkit/docs/quickstarts/slurm-cluster>

```
srun -N 3 hostname
```

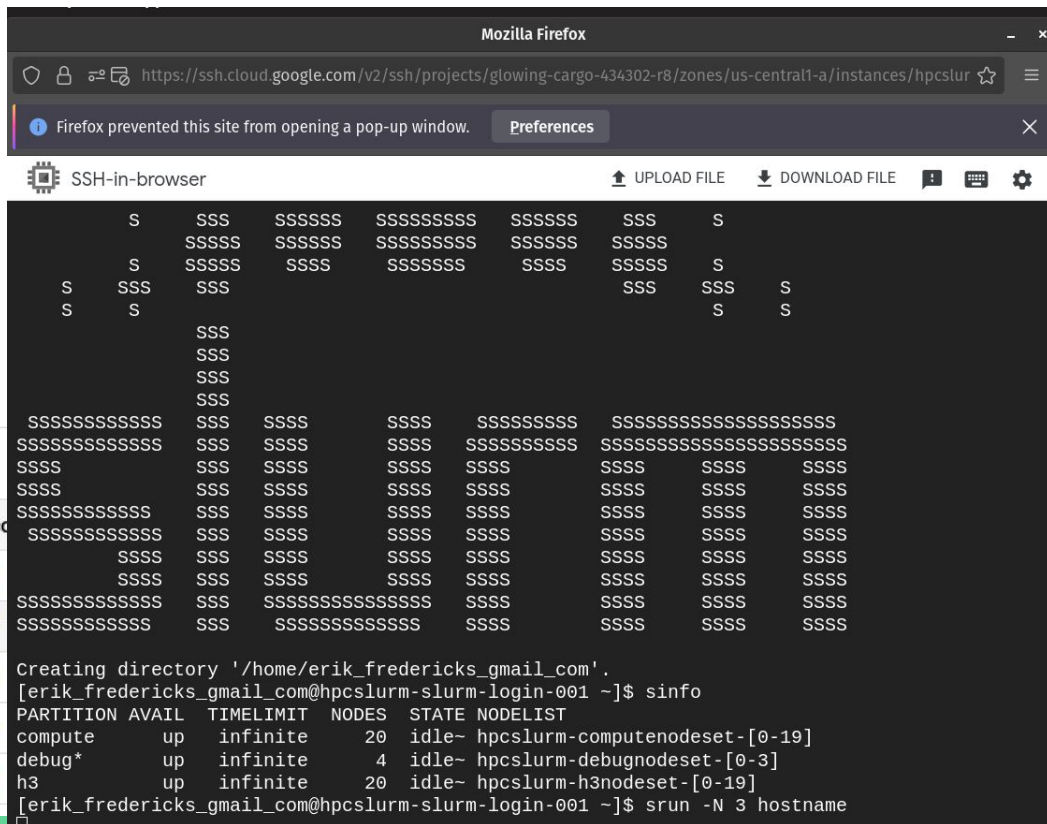
This command creates three compute nodes for your HPC cluster. This may take a minute while Slurm auto-scales to create the three nodes.

When the job finishes you should see an output similar to:

```
$ srun -N 3 hostname
hpcslurm-debug-ghpc-0
hpcslurm-debug-ghpc-1
hpcslurm-debug-ghpc-2
```

VM instances

<input type="checkbox"/>	Status	Name ↑	Zone	Rec...
<input type="checkbox"/>	✓	hpcslurm-controller	us-central1-a	
<input type="checkbox"/>	'	hpcslurm-debugnodeset-0	us-central1-a	'
<input type="checkbox"/>	'	hpcslurm-debugnodeset-1	us-central1-a	'
<input type="checkbox"/>	'	hpcslurm-debugnodeset-2	us-central1-a	'
<input type="checkbox"/>	✓	hpcslurm-slurm-login-001	us-central1-a	



PaaS

Before we didn't worry about the bare-metal hardware

- Now we don't worry about that **plus** the operating system/environment!
- i.e., a *serverless* approach

i.e., we just want to deploy a Python app and don't care about the server itself

Google Cloud - App Engine / Cloud Functions

AWS - App Runner / Lambda

Azure - Azure PaaS / Azure Functions

(there are more - just these are common)

PaaS Demo!

This one's pretty straightforward - basic website with App Engine

<https://cloud.google.com/appengine/docs/standard/hosting-a-static-website>

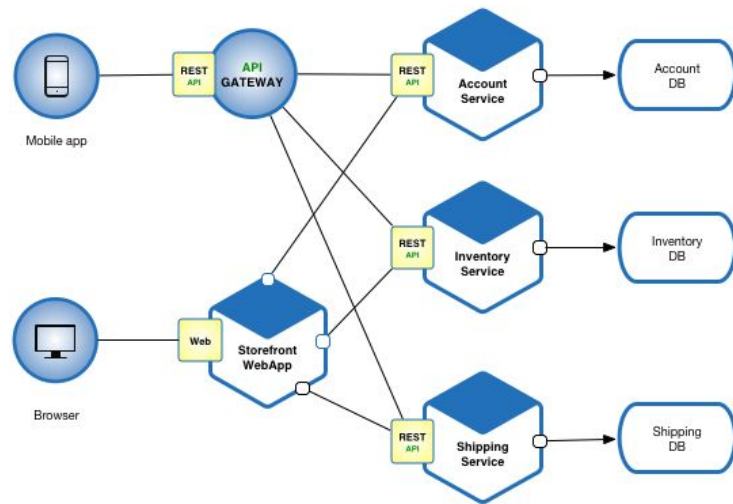
Buuut,

- Gets us to YAML (get ready for pain and suffering)

Functions as a Service (FaaS)

Less common in marketing materials, but here:

- Functions deployed remotely without a care for the app *as a whole*
- i.e., I am calling a function as an API
 - Part of the serverless paradigm, so still PaaS
- Lambda / Cloud Functions / Azure Functions also apply here



What are some sample use cases?

We'll be going into greater detail on these later

SaaS

Don't care about the hardware

- Or the operating system
- Or the installation!

Think: Google Docs, Office 365, etc.



And now ... MaaS

AI and machine learning products

Try **Gemini 1.5 models**, the latest and most advanced multimodal models in Vertex AI. See what you can build with up to a 2M token context window, starting as low as \$0.0001.

[Try it in console](#)[Contact sales](#)

Summarize large documents with generative AI

Deploy a preconfigured solution that uses generative AI to quickly extract text and summarize large documents.



Deploy an AI/ML image processing pipeline

Launch a preconfigured, interactive solution that uses pre-trained machine learning models to analyze images and generate image annotations.



Create a chat app using retrieval-augmented generation (RAG)

Deploy a preconfigured solution with a chat-based experience that provides questions and answers based on embeddings stored as vectors.



Ok, so we have the various *aaS

Things **you** need to take away from this:

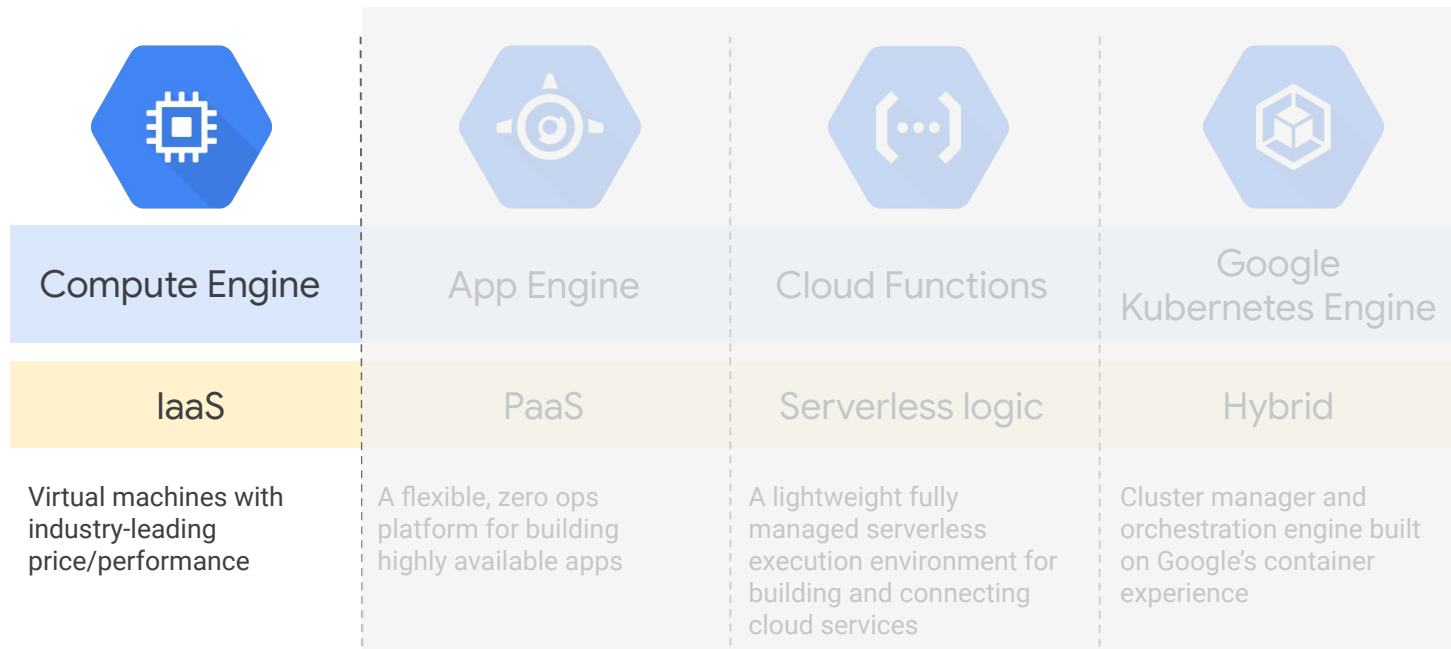
- 1) If you have a need, there's probably a solution out there
- 2) Know **what** technology to use for **which** problem

FOR EXAMPLE - WHAT TECH CAN YOU USE?

- 1) I need you to build a website for our company
- 2) I need you to create a backend service so our users can purchase products from the online store
- 3) I need you to create a chatbot because [expletive deleted] every website needs a chatbot now

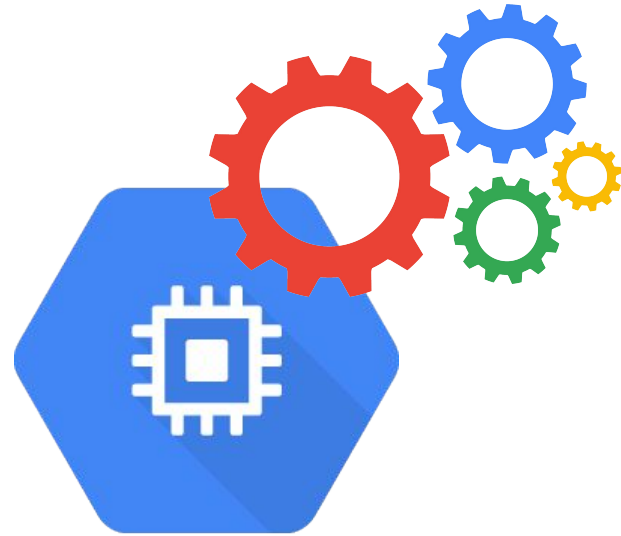
Let's take a look at IaaS for a moment

Where Compute Engine fits within Google Cloud



Compute Engine is an infrastructure-centric solution

- Type of IaaS
- Scalable, high-performance VMs
- Run any computing workload
- Predefined or custom machine types
- Windows or Linux
- No upfront investment required



Create VMs that are right for your workloads

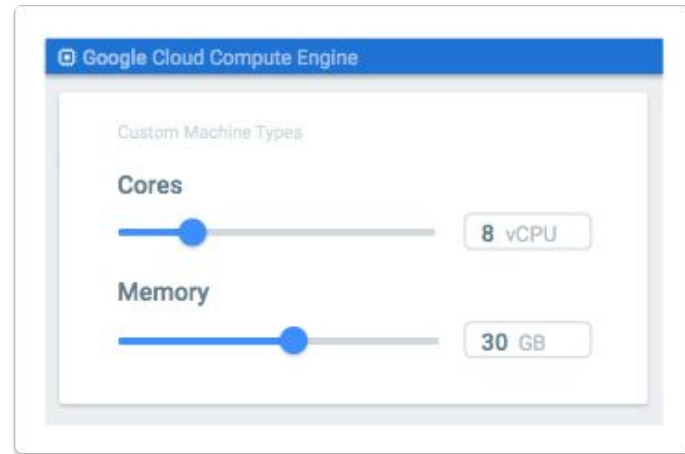
Machine type options to consider:

- Higher proportion of memory to CPU
- Higher proportion of CPU to memory
- Blend of both

Select from predefined VM configurations:

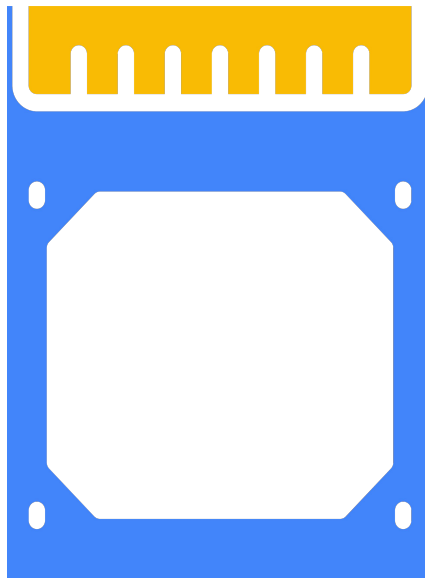
- General-purpose
- Memory-optimized
- Compute-optimized

Customize your own configuration



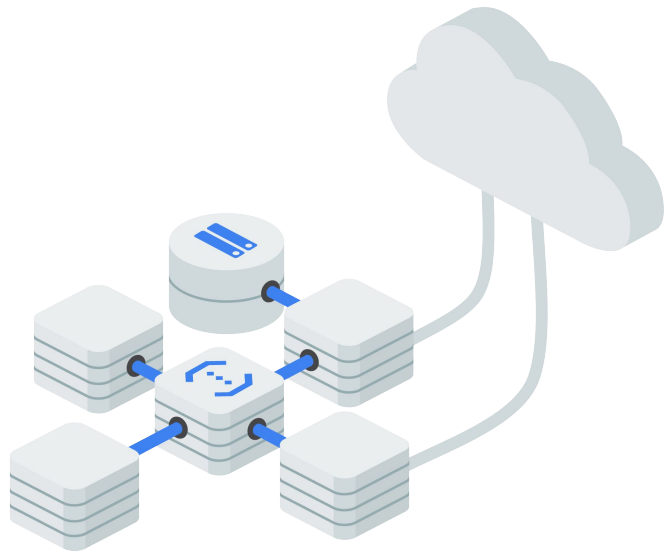
Building virtual disks

- Network storage can be attached to VMs as persistent disks (PDs).
- PDs are cost-effective, durable, and offer good performance.
- Local SSDs provide higher performance with lower latency, but exist only for the lifetime of a specific instance.
- Standard PD throughput performance and IOPS increases linearly with the size of the disk until it reaches set per-instance limits.
- SSD PD IOPS performance depends on the number of vCPUs in the instance in addition to disk size.



Networks connect Compute Engine instances to each other and to the internet

- Inbound/outbound firewall rules
- Create static routes
- Scale and distribute applications using Cloud Load Balancing
- Global and multi-regional subnetwork



Compute Engine pricing

Google Cloud Platform Pricing Calculator

Prices are up to date. Last update: 15-August-2019



Search for a product you are interested in.

Instances

Number of instances *

5

What are these instances for?

Databases

Operating System / Software

Free: Debian, CentOS, CoreOS, Ubuntu, or other User Provided OS

Machine Class

Preemptible

Machine Family

General purpose

Machine Generation

Second (latest)

Machine type

n2-standard-2 (vCPUs: 2, RAM: 8GB)



Estimate

Compute Engine

5 x Databases



3,650 total hours per month

VM class: preemptible

Instance type: n2-standard-2

Region: Iowa

Total available local SSD space 1x375 GB

Estimated Component Cost: USD 175.77 per 1 month

Total Estimated Cost: USD 175.77 per 1 month

Estimate Currency

USD - US Dollars

EMAIL ESTIMATE

SAVE ESTIMATE

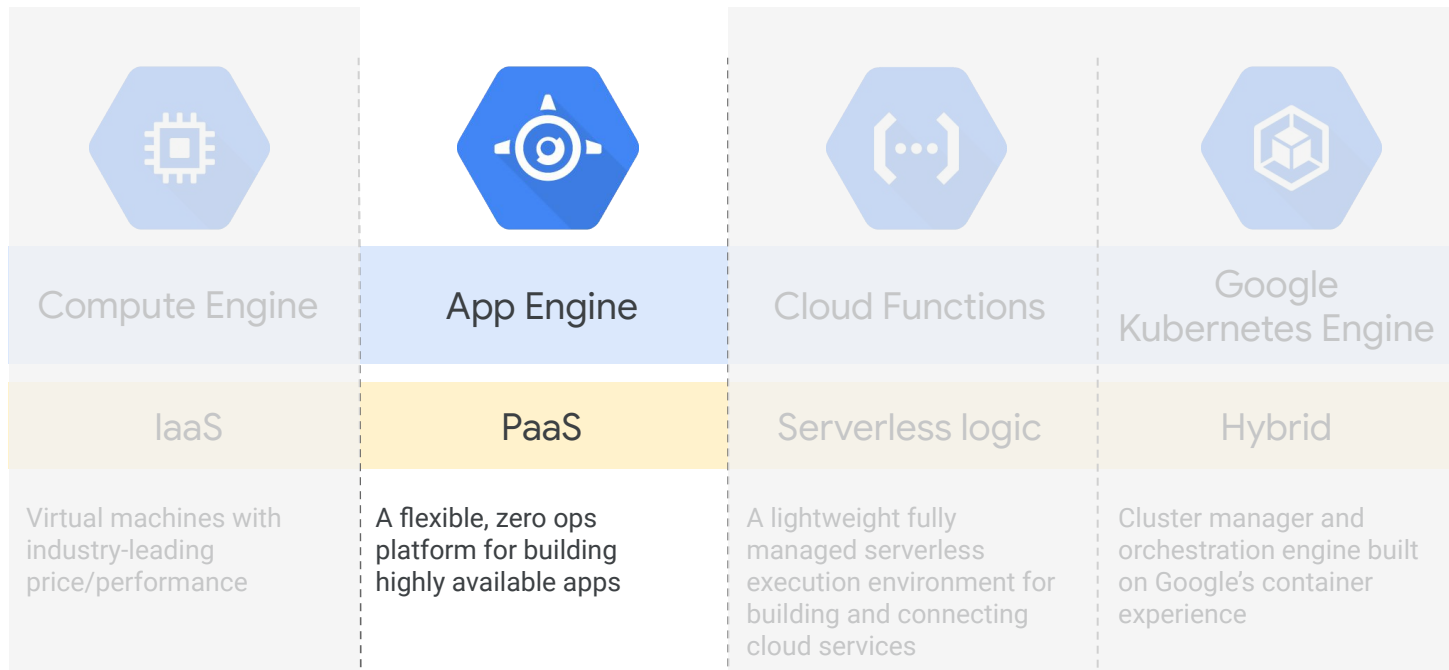
<https://cloud.google.com/products/calculator/>



and now PaaS

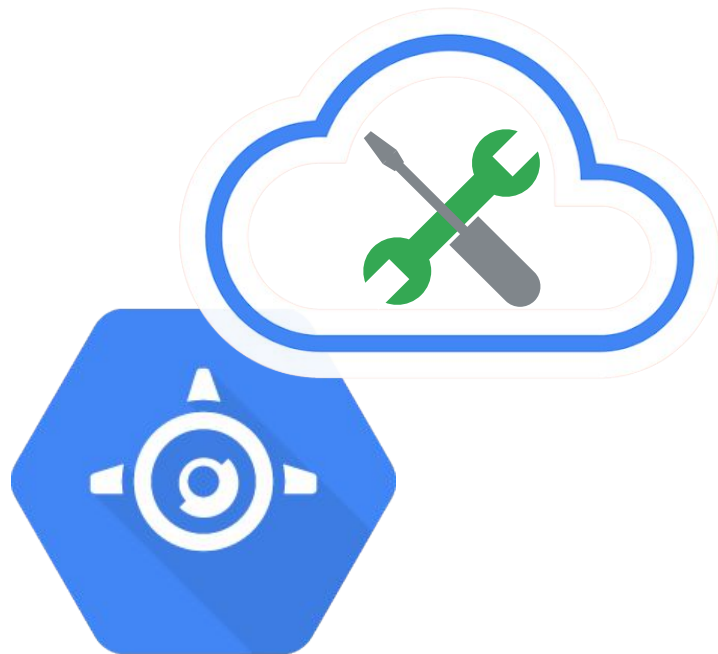


Where App Engine fits within Google Cloud



App Engine is a platform-centric solution

- Type of PaaS
- No need to buy, build, or operate hardware/infrastructure
- No managing servers or configuring deployments
- Focus on app development instead of operations
- Use a range of languages and tools
- Automatic scaling



App Engine offers two different environments

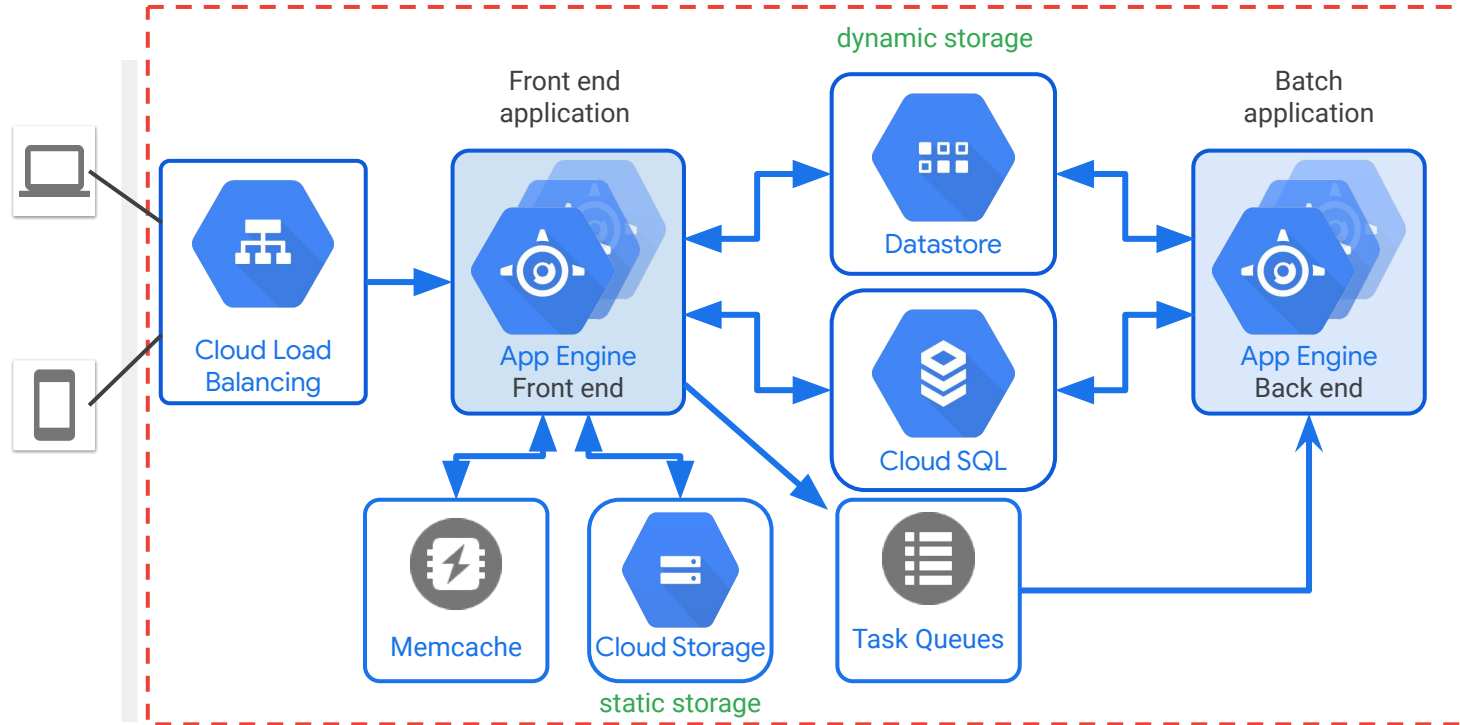
Standard environment

- Fully-managed
- Scale to zero
- Specific versions of supported languages
- Changes/configuration limited

Flexible environment

- Docker container support
- VMs exposed
- Any language in your container
- More options for infrastructure customization and configuration for performance

An App Engine architecture example



App Engine addresses the key needs of developers



Multiple storage options



Automatic scaling



Load balancing



App versioning



Monitoring and logging



Security

SaaS

<https://youtu.be/WJmP4Lqopps>

- (Important bit? What tools/tech to use to build the "thing" you want)

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

- Tends to be 'install trusted plugin from Marketplace'

Another demo!

<https://developers.google.com/codelabs/gcp-marketplace-saas#0>

<appears broken>

Let's do a PHP app! Different than what we've done so far

<https://cloud.google.com/appengine/docs/standard/php-gen2/building-app>