

GCP

Professional Cloud Architect



Agenda

- 1. Introduction**
- 2. Business Logic Design**
- 3. Data Layer Design**
- 4. Network Layer Design**
- 5. Design for security**

Introduction

What is a Cloud Architect?

- A cloud Architect enables organizations to leverage...Cloud technologies.

Through an understanding of cloud architecture and ...technology, this individual designs, develops, and manages robust, secure, scalable, highly available, and dynamic solutions to drive business objectives.

- The Cloud architect should be proficient in all aspects of enterprise cloud strategy, solution design, and architectural best practices.
- The cloud Architect should also be experienced in software development methodologies and approaches including multi-tiered distributed applications which span multi-cloud or hybrid environments.

Key abilities

The Cloud Architect [should have the] ability to:

- Design and plan a cloud solution architecture**
- Manage and provision the cloud solution infrastructure**
- Design for security and compliance**
- Analyze and optimize technical and business processes**
- Manage implementations of cloud architecture**
- Ensure solution and operations reliability**

What are “Systems”?

- The “Things” we architect and build with cloud tools
- Always about data flows:
 - Moving information
 - Remembering information
 - Processing information
- Driven by a purpose
 - To achieve something
 - To solve a problem

Not just "Architecting Great Systems"

Not just great *systems*, when they're in use...

Nor just great *architectures*, when they're complete...

But also great *architecting*, while it happens...

And great *building* of those architectures.

The *process* matters.

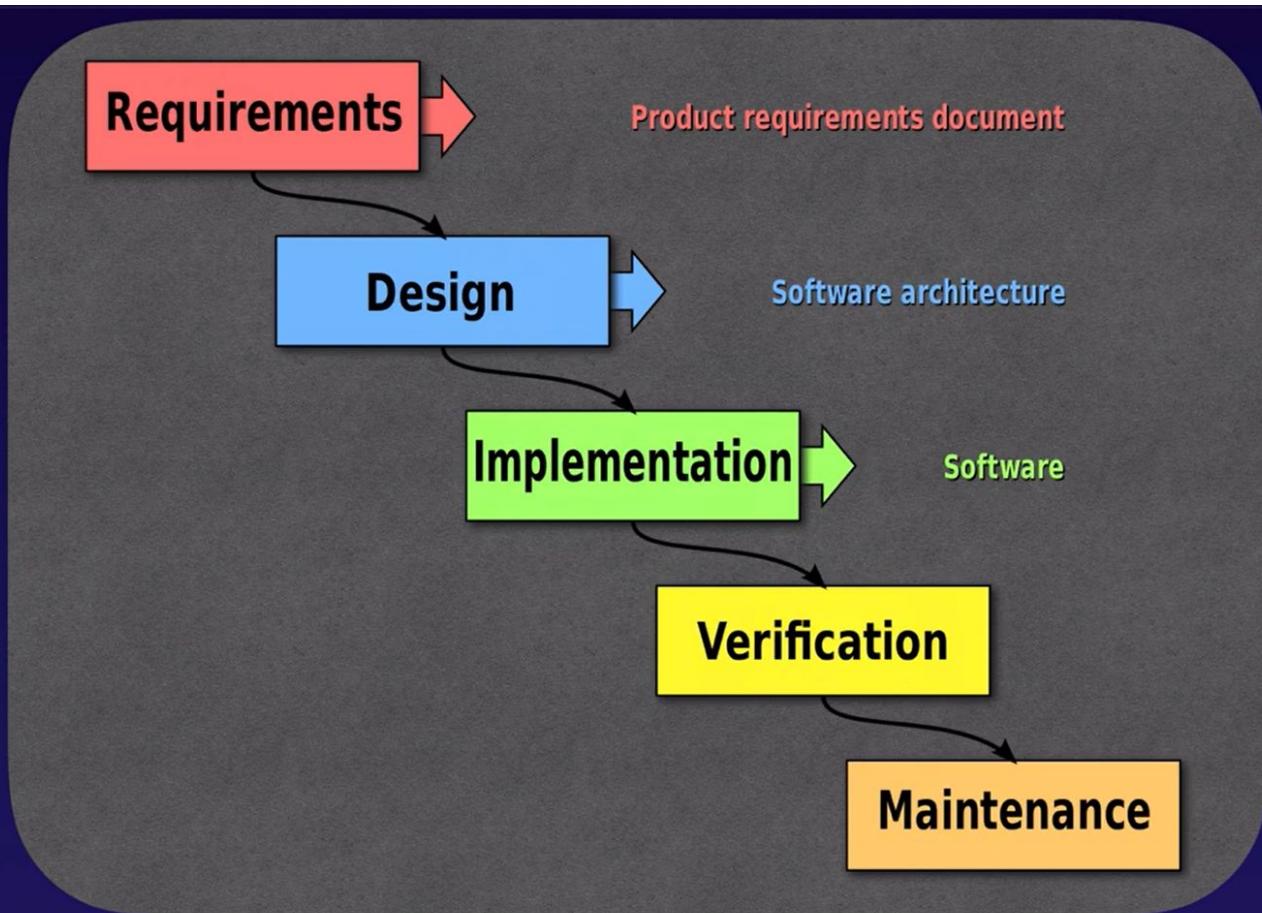
A lot.

Positionally



**Cloud Architect
is a
Leadership
Position**

Waterfall



Agile

Manifesto for Agile Software Development

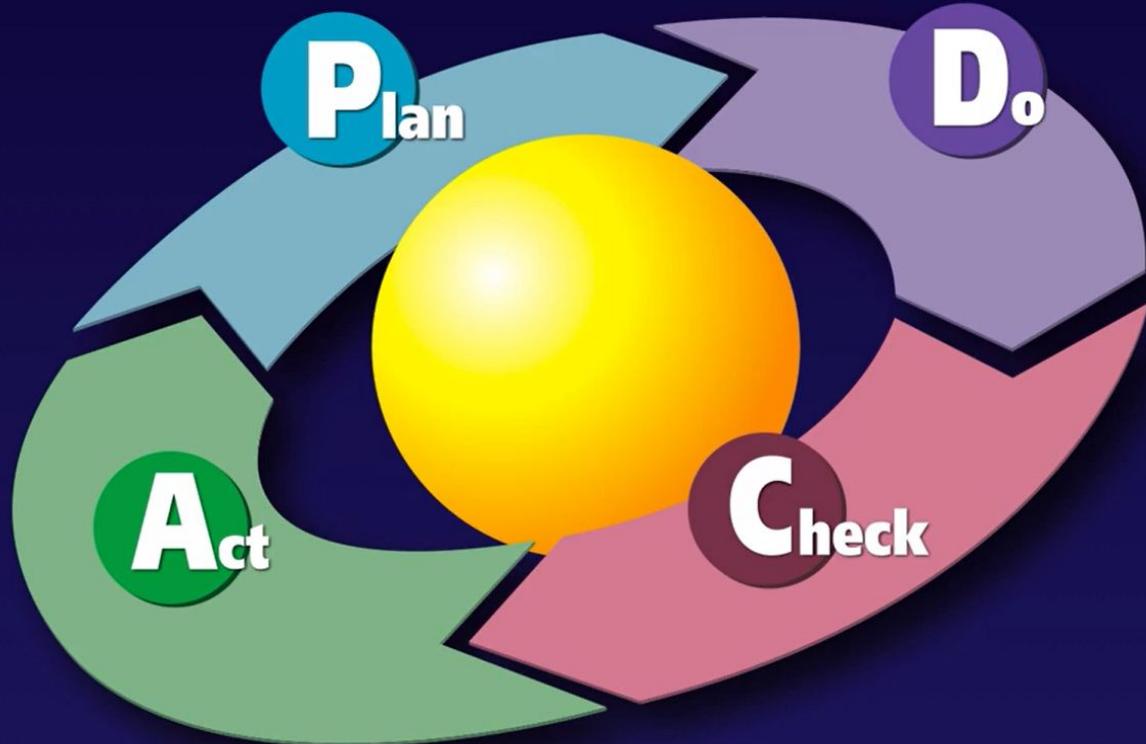
We are uncovering better ways of developing software by doing it and helping others do it.
Through this work we have come to value:

Individuals and interactions over processes and tools
Working software over comprehensive documentation
Customer collaboration over contract negotiation
Responding to change over following a plan

That is, while there is value in the items on the right, we value the items on the left more.



Iteration: Deming Cycle/PDCA



Strategic Stuff

How to Plan:

- 1. Understand your environment**
- 2. Determine where you are now**
- 3. Figure out where you want to go**
- 4. Decide how you'll get there**

Doing the Architecting Thing

Foundation	How does this stuff work?
Environment	What affects us?
Current	Where are we?
Target	Where do we want to be?
Plan	How will we get from here to there?
Execute	Let's go!
Evaluate	Are we getting there?
Adjust	What should we change?

Business Logic Design

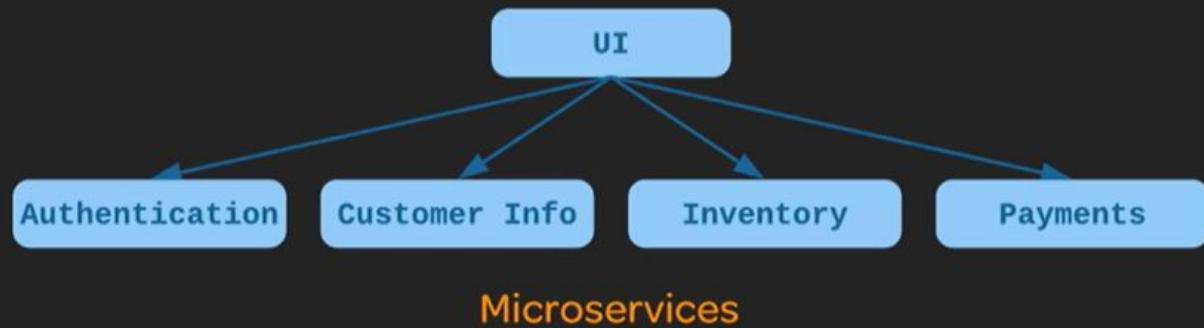
What are Microservices?

- Microservices: A microservice architecture breaks an application up into a collection of small, loosely-coupled services
- Traditionally, apps used a monolithic architecture. In a monolithic architecture, all features and services are part of one large application
- Microservices are small: each microservice implements only a small piece of an application's overall functionality
- Microservices are loosely coupled: Different microservices interact with each other using stable and well-defined APIs. This means that they are independent of one another

Microservices vs Monolith



Monolith



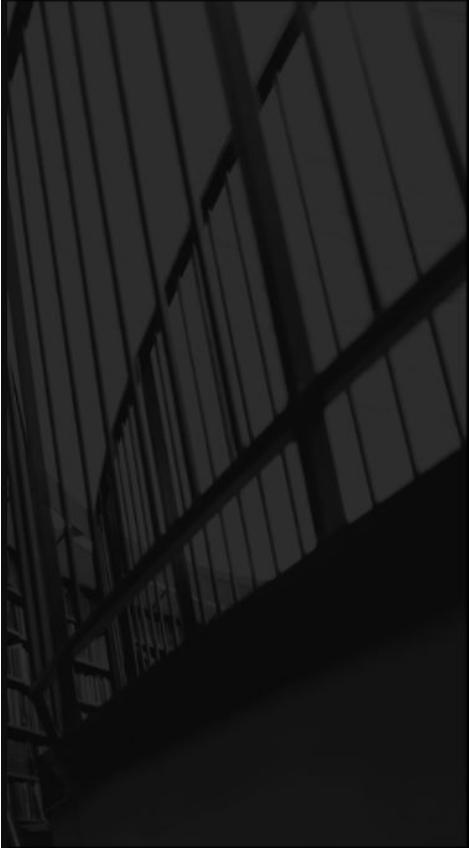
What do microservices look like?

- There are many different ways to structure and organize a microservice architecture
- For example, a pet shop application might have:
 - A pet inventory service
 - A customer details service
 - An authentication service
 - A pet adoption request service
 - A payment processing service
- Each of these is its own codebase and a separate running process (or processes). They can all be built, deployed, and scaled separately

Why use Microservices?

- Modularity – Microservices encourage modularity. In monolithic apps, individual pieces become tightly coupled, and complexity grows. Eventually, it's very hard to change anything without breaking something
- Technological flexibility – You don't need to use the same languages and technologies for every part of the app. You can use the best tool for each job
- Optimized scalability – You can scale individual parts of the app based upon resource usage and load. With a monolith, you have to scale up the entire application, even if only one aspect of the service actually needs to be scaled
- Microservices aren't always the best choice. For smaller, simpler apps a monolith might be easier to manage

The Twelve-Factor App:

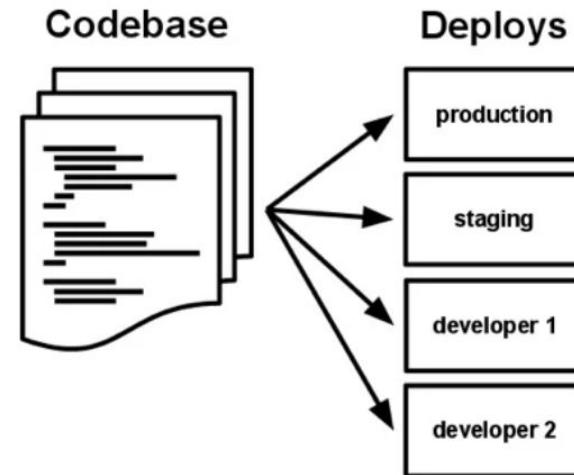


1. Codebase
2. Dependencies
3. Config
4. Backing Services
5. Build, Run, Release
6. Stateless Processes
7. Port Binding
8. Concurrency
9. Disposability
10. Dev-Prod Parity
11. Logs
12. Admin Processes

The twelve factors - Codebase

CODEBASE

- codebase = repo
- **one** repo - **many** deploys
- app **!=** many repos
- many repos = distributed system



The twelve factors - Dependencies

DEPENDENCIES

Explicitly declare and isolate dependencies

1. Declare **dependencies** in a **manifest**
2. Use **isolation** tools
3. **Specific** versions are important
4. Avoid **shelling** to unbundled system tools.

The twelve factors - Configuration

- Config is the **specific** information required to host a deployment of your codebase
- Does not include things like routes etc.
- Mainly database **credentials, paths, resource urls** etc.

The twelve factors - Configuration

Examples:

Dependency Manifest = **Gemfile**

Isolation tools = **bundle exec**

```
gem "redis-rails", "~> 3.2.3"
```

The twelve factors - Configuration

Store config in the environment.

- Keep your config **outside** the app
- No config in git
- Open source test

The twelve factors - Configuration

Use **environment vars**

Can you make your repo open source
today?

ENV['DATABASE_URL']

```
production:  
  adapter: mysql2  
  database: <%= ENV['DB_ENV_MYSQL_DATABASE'] %>  
  username: <%= ENV['DB_ENV_MYSQL_USER'] %>  
  password: <%= ENV['DB_ENV_MYSQL_PASS'] %>  
  port: <%= URI(ENV['DB_PORT']).port %>  
  host: <%= URI(ENV['DB_PORT']).host %>
```

The twelve factors - Backing services

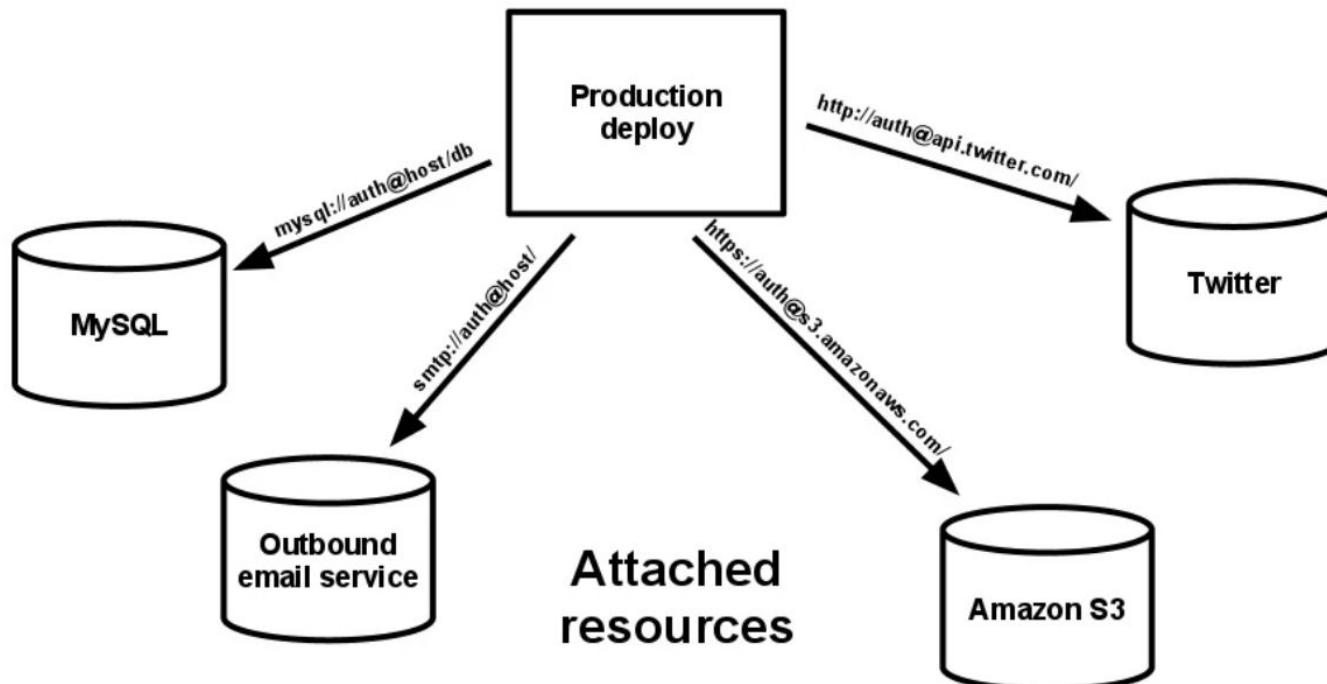
Treat backing services as attached resources

What's a backing service?

- **Datastore**
- **SMTP**
- **Caching systems**
- **Amazon S3**

Make no distinction between local and third party services

The twelve factors - Backing services



BUILD = codebase + dependencies + assets

RELEASE = **BUILD** + config

RUN = run process against **RELEASE**

- Strict **separation** between stages
- **Cannot change** code at **runtime**
- **Rollback** = just use the last release instead.
- Release has unique release **ID**

The twelve factors - Build, release, run

- Does the **running** of the **release**
- Is **stateless**
- **Shares nothing** with other processes
- Uses **single transaction** only caches
- Session **db storage** ...over sticky sessions
- Asset **pre-compilation** ...over runtime calculation

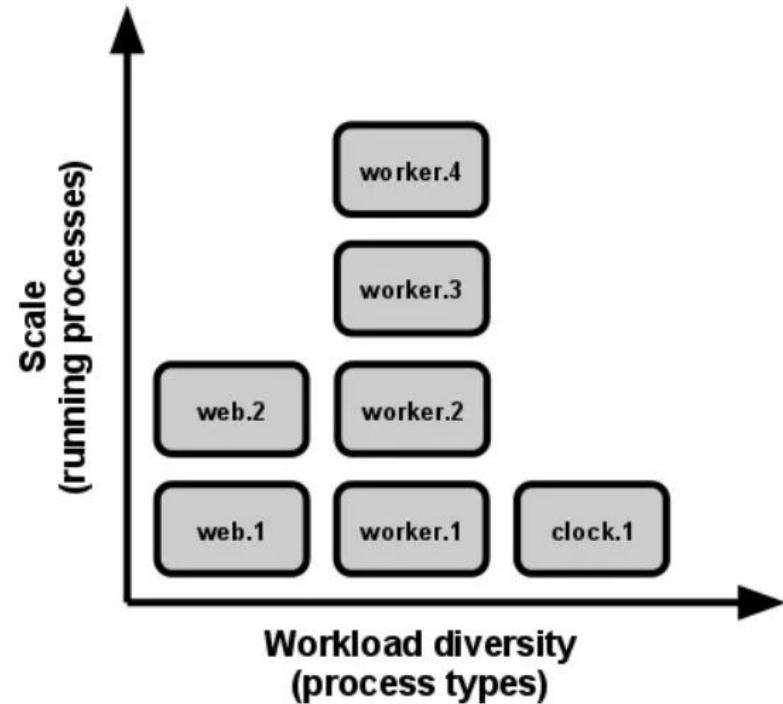
The twelve factors - Processes

The twelve factors - Port binding

- **Scale out** via the **process model**
- Processes are **first class** citizens.
- Assign work to a process **type** (web, db, worker etc.)
- Can then handle **own multiplexing**
- Processes **don't daemonize** or write PID files.
- Instead, use system **process manager**

The twelve factors - Concurrency

Scale out by running multiple processes of different types



Processes are **disposable**

- Start **quickly**
- Shut down **gracefully**
- Be robust against **sudden death**

Example: Worker Process

- Return the current job to the job queue
- All jobs are **reentrant**
- Or at least **idempotent**

Gaps exist between development and production:

- **Time** (days to push)
- **Personnel** (dev vs ops)
- **Tools**

Need to **shorten the gaps!**

- CI & deploy **ASAP** after writing code
- Get Developers **involved** in Operations
- Environments should be as **similar** as possible:
 - Eg. **Resist the urge** to use **different** backing services between development and production

Within your app treat logs as **event streams**

- **Don't** route or store **logs in files**
- **Stream** to stout instead and be captured and handled by the **environment**

The twelve factors - Admin processes

One-off admin tasks include:

- Database **migrations**
- **Console**
- One-time **scripts**

The twelve factors - Admin processes

- Run as **separate** process
- Run against the **same release**
- **Admin code** ships with **app code**

Mapping computational needs for Google Cloud Platform processing services

Data Layer Design

Data flow foundations

Network
Compute
Storage

Moving
Processing
Remembering



Data Flows - Processing:

Data Flows - Processing: Compute Engine (GCE)



Compute
Engine

Zonal Regional Multi-Regional Global



- **Fast-booting Virtual Machines (VMs) you can rent, on demand**
- **Infrastructure as a Service (IaaS)**
- **Pick set machine type—standard, highmem, highcpu—or custom CPU/RAM**
- **Pay by the second (60 second min.) for CPUs, RAM**
- **Automatically cheaper if you keep running it (“sustained use discount”)**
- **Even cheaper for “preemptible” or long-term use commitment in a region**
- **Can add GPUs and paid OSes for extra cost**
- **Live Migration: Google seamlessly moves instance across hosts, as needed**

Data Flows - Processing: Kubernetes Engine (GKE)



Kubernetes
Engine

Zonal Regional Multi-Regional Global



- Managed Kubernetes cluster for running Docker containers (with autoscaling)
- Used to be called “Google Container Engine” (but still GKE) until Nov, 2017
- Kubernetes DNS on by default for service discovery
- No IAM integration (unlike AWS’s ECS)
- Integrates with Persistent Disk for storage
- Pay for underlying GCE instances
 - Production cluster should have 3+ nodes
- No GKE management fee, no matter how many nodes in cluster, as of Nov 2017

Data Flows - Processing: App Engine (GAE)



App
Engine

Zonal Regional Multi-Regional Global



- Platform as a Service (PaaS) that takes your code and runs it



Elastic
Beanstalk



Data Flows - Processing: App Engine (GAE)



App
Engine

Zonal Regional Multi-Regional Global



- Platform as a Service (PaaS) that takes your code and runs it
- Much more than just compute — Integrates storage, queues, NoSQL, ...
- Flex mode ("App Engine Flex") can run any container & access VPC
- Auto-scales based on load
 - Standard (non-Flex) mode can turn off last instance when no traffic
- Effectively pay for underlying GCE instances and other services

Data Flows - Processing: Cloud Functions (GCF)



Cloud
Functions

Zonal Regional Multi-Regional Global



- Runs code in response to an event — Node.js, Python, Java, Go
- Functions as a Service (FaaS), “Serverless”
- Pay for CPU and RAM assigned to function, per 100ms (min. 100ms)
- Each function automatically gets an HTTP endpoint
- Can be triggered by GCS objects, Pub/Sub messages, etc.
- Massively scalable (horizontally) — Runs many copies when needed
- Often used for chatbots, message processors, IoT, automation, etc.

Choosing compute: Flowchart

Are you a mobile/HTML5 developer?

Y



Firebase

N

Are you developing an event-driven app?

Y



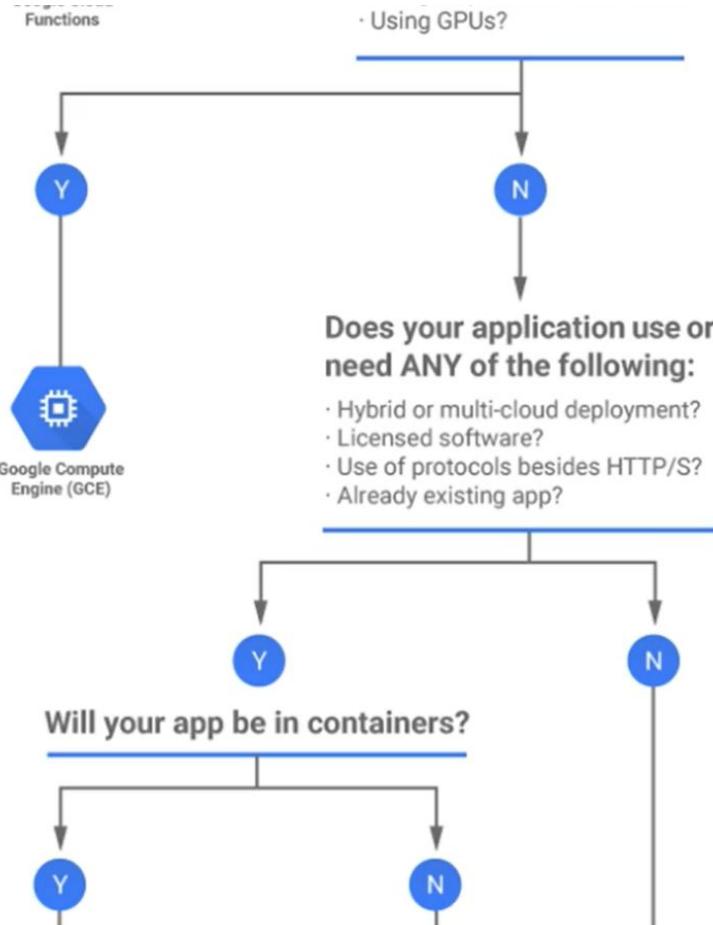
Google Cloud
Functions

N

Is your application
ANY of the following:

- Using a specific OS or kernel?
- Using GPUs?

Choosing compute: Flowchart



Choosing compute: Flowchart



Google Compute
Engine (GCE)

Does your application use or
need ANY of the following:

- Hybrid or multi-cloud deployment?
- Licensed software?
- Use of protocols besides HTTP/S?
- Already existing app?

Y

Will your app be in containers?

N

Y

Do you plan to use Kubernetes
as your orchestrator?

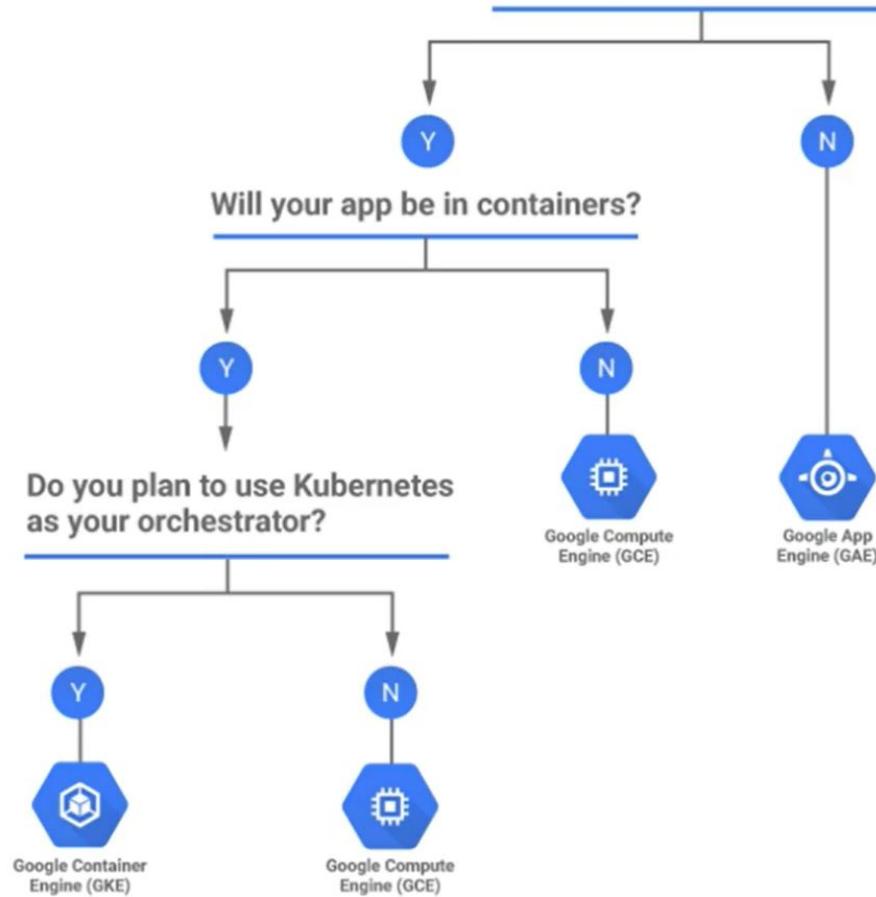
N

Google Compute
Engine (GCE)

Google App
Engine (GAE)

N

Choosing compute: Flowchart



Data Flows - Processing: Autoscaling Compute Comparison

Autoscaling	Minimum	Min for HA	Mechanism	Speed
Compute Engine	1 instance	2 instances	MIG alone / MIG + CLB	Moderate / Better
	1 pod	2-pod deployment	Load-based internal K8s	Decent
GKE Node Pool	1 instance	3 instances	K8s pod back-pressure	Slow
	1 container	2 containers	Request-based	"Gradually"
App Engine Standard	Zero	Zero	Request-based	"Sudden and extreme spikes of traffic" OK
Cloud Functions	Zero	Zero	Request-based	Very Fast

Data Flows - Processing: Machine Learning and AI



Cloud
ML Engine

Zonal Regional Multi-Regional Global



- Massively scalable managed service for training ML models & making predictions



Amazon
SageMaker



TensorFlow



Apache
MXNet

Data Flows - Processing: Machine Learning and AI



Cloud
ML Engine

Zonal Regional Multi-Regional Global



- Massively scalable managed service for training ML models & making predictions
- Enables apps/devs to use TensorFlow on datasets of any size; endless use cases
- Integrates: GCS/BQ (storage), Cloud Datalab (dev), Cloud Dataflow (preprocessing)
- Supports online & batch predictions, prioritizing latency (online) & job time (batch)
- Or download models & make predictions anywhere: desktop, mobile, own servers
- HyperTune automatically tunes model hyperparameters to avoid manual tweaking
- Training: Pay per hour depending on chosen cluster capabilities (ML training units)
- Prediction: Pay per provisioned node-hour plus by prediction request volume made

Data Flows - Processing: Machine Learning and AI



Cloud
Vision API

Zonal Regional Multi-Regional Global



- Classifies images into categories, detects objects/faces, & finds/reads printed text
- Pre-trained ML model to analyze images and discover their contents
- Classifies into thousands of categories (e.g., "sailboat", "lion", "Eiffel Tower")
- Upload images or point to ones stored in GCS
- Pay per image, based on detection features requested
 - Higher price for OCR of full documents and finding similar images on the web
 - Some features are priced together: Labels + SafeSearch, ImgProps + Cropping
 - Other features priced individually: Text, Faces, Landmarks, Logos

Data Flows - Processing: Machine Learning and AI



Cloud
Speech API

Zonal Regional Multi-Regional Global



- Automatic Speech Recognition (ASR) to turn spoken word audio files into text
- Pre-trained ML model for recognizing speech in 110+ languages/variants
- Accepts pre-recorded or real-time audio, & can stream results back in real-time
- Enables voice command-and-control and transcribing user microphone dictations
- Handles noisy source audio
- Optionally filters inappropriate content in some languages
- Accepts contextual hints: words and names that will likely be spoken
- Pay per 15 seconds of audio processed

Data Flows - Processing: Machine Learning and AI



Cloud Natural
Language API

Zonal Regional Multi-Regional Global



- Analyzes text for sentiment, intent, & content classification, and extracts info
- Pre-trained ML model for understanding what text means, so you can act on it
- Excellent with Speech API (audio), Vision API (OCR), & Translation API (or built-ins)
- Syntax analysis extracts tokens/sentences, parts of speech & dependency trees
- Entity analysis finds people, places, things, etc., labels them, & links to Wikipedia
- Analysis for sentiment (overall) and entity sentiment detect +/- feelings & strength
- Content classification puts each document into one of 700+ predefined categories
- Charged per request of 1000 characters, depending on analysis types requested

Data Flows - Processing: Machine Learning and AI



Cloud
Translation API

Zonal Regional Multi-Regional Global



- Translate text among 100+ languages; optionally auto-detects source language
- Pre-trained ML model for recognizing and translating semantics, not just syntax
- Can let people support multi-regional clients in non-native languages, 2-way
- Combine with Speech, Vision, & Natural Language APIs for powerful workflows
- Send plain text or HTML and receive translation in kind
- Pay per character processed for translation
- Also pay per character for language auto-detection

Data Flows - Processing: Machine Learning and AI



Dialogflow

Zonal Regional Multi-Regional Global



- Build conversational interfaces for websites, mobile apps, messaging, IoT devices
- Pre-trained ML model and service for accepting, parsing, lexing input & responding
- Enables useful chatbots and other natural user interactions with your custom code
- Train it to identify custom entity types by providing a small dataset of examples
- Or choose from 30+ pre-built agents (e.g. car, currency, dates) as starting template
- Supports many different languages and platforms/devices
- Free plan has unlimited text interactions and capped voice interactions
- Paid plan is unlimited but charges per request: more for voice, less for text

Data Flows - Processing: Machine Learning and AI



Cloud Video
Intelligence API

Zonal Regional Multi-Regional Global



- Annotates videos in GCS (or directly uploaded) with info about what they contain
- Pre-trained ML model for video scene analysis and subject identification
- Enables you to search a video catalog the same way you search text documents
- “Specify a region where processing will take place (for regulatory compliance)”
- *Label Detection:* Detect entities within the video, such as "dog", "flower" or "car"
- *Shot Change Detection:* Detect scene changes within the video
- *SafeSearch Detection:* Detect adult content within the video
- Pay per minute of video processed, depending on requested detection modes

Data Flows - Processing: Machine Learning and AI



Cloud Job
Discovery

Zonal Regional Multi-Regional Global



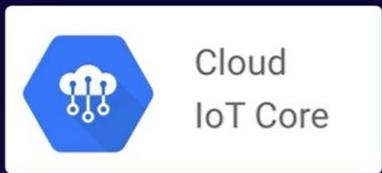
- Helps career sites, company job boards, etc. to improve engagement & conversion
- Pre-trained ML model to help job seekers search job posting databases
- Most job sites rely on keyword search to retrieve content which often omits relevant jobs and overwhelms the job seeker with irrelevant jobs. For example, a keyword search with any spelling error returns 0 results, and a keyword search for "dental assistant" returns any "assistant" role that offers dental benefits.'
- Integrates with many job/hiring systems
- Lots of features, such as commute distance and recognizing abbreviations/jargon
- "Show me jobs with a 30 minute commute on public transportation from my home"

Data Flows - Processing: Big Data Lifecycle

Ingest	Store	Process & Analyze	Explore & Visualize
 App Engine	 Cloud Storage	 Cloud Dataflow	 Cloud Datalab
 Compute Engine	 Cloud SQL	 Cloud Dataproc	 Google Data Studio
 Container Engine	 Cloud Datastore	 BigQuery	 Google Sheets
 Cloud Pub/Sub	 Cloud Bigtable	 Cloud ML	
 Stackdriver Logging	 BigQuery	 Cloud Vision API	
 Cloud Transfer Service		 Cloud Speech API	
		 Translate API	
		 Cloud Natural Lang API	

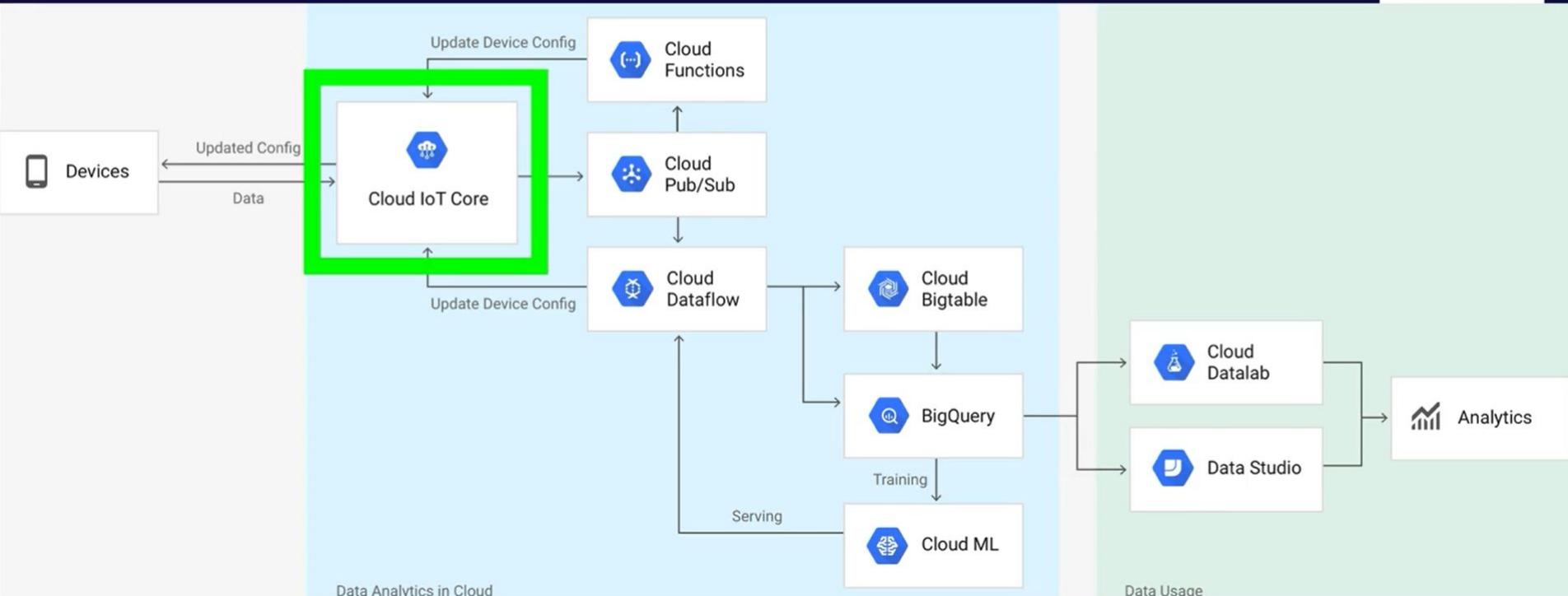


Data Flows - Processing: Big Data and IoT



Cloud
IoT Core

Zonal Regional Multi-Regional Global



Data Flows - Processing: Big Data and IoT



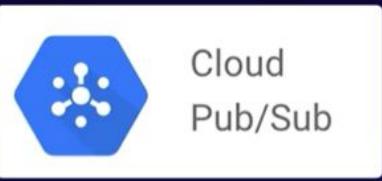
Cloud
IoT Core

Zonal Regional Multi-Regional Global



- Fully-managed service to connect, manage, and ingest data from devices globally
- Device Manager handles device identity, authentication, config, & control
- Protocol Bridge publishes incoming telemetry to Cloud Pub/Sub for processing
- Connect securely using IoT industry-standard MQTT or HTTPS protocols
- CA signed certificates can be used to verify device ownership on first connect
- Two-way device communication enables configuration & firmware updates
- Device shadows enable querying & making control changes while devices offline

Data Flows - Processing: Big Data and IoT

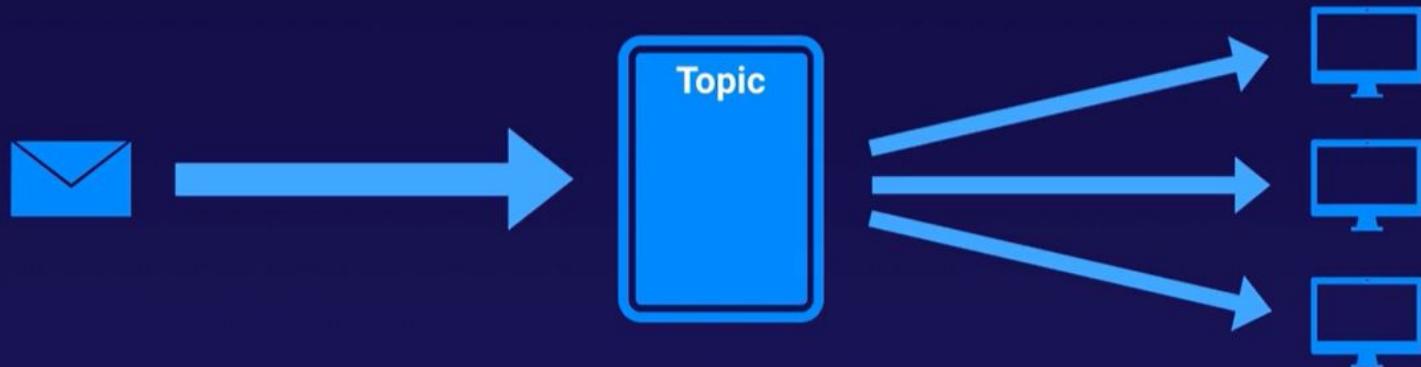


Cloud
Pub/Sub

Zonal Regional Multi-Regional Global



- Infinitely-scalable at-least-once messaging for ingestion, decoupling, etc.



Data Flows - Processing: Big Data and IoT



Cloud
Pub/Sub

Zonal Regional Multi-Regional Global



- Infinitely-scalable at-least-once messaging for ingestion, decoupling, etc.



Amazon
SNS



Amazon
SQS



RabbitMQ



Amazon
Kinesis



Apache
Kafka

Data Flows - Processing: Big Data and IoT



Cloud
Pub/Sub

Zonal Regional Multi-Regional Global



- Infinitely-scalable at-least-once messaging for ingestion, decoupling, etc.
- “Global by default: Publish... and consume from anywhere, with consistent latency.”
- Messages can be up to 10MB and undelivered ones stored for 7 days—but no DLQ
- Push mode delivers to HTTPS endpoint & succeeds on HTTP success status code
 - “Slow-start” algorithm ramps up on success and backs off & retries, on failures
- Pull mode delivers messages to requesting clients and waits for ACK to delete
 - Lets clients set rate of consumption, and supports batching and long-polling
- Pay for data volume
 - Min 1KB per publish/push/pull request (not by message)

Data Flows - Processing: Big Data and IoT



Cloud
Dataprep

Zonal Regional Multi-Regional Global



- Visually explore, clean, and prepare data for analysis without running servers
- “Data Wrangling” (i.e. “ad-hoc ETL”) for business analysts, not IT pros
 - Who might otherwise spend 80% of their time cleaning data
- Managed version of Trifacta Wrangler—and managed by Trifacta, not Google
- Source data from GCS, BQ, or file upload—formatted in CSV, JSON, or relational
- Automatically detects schemas, datatypes, possible joins, and various anomalies
- Pay for underlying Dataflow job, plus management overhead charge



Data Flows - Processing: Big Data and IoT



Cloud
Dataproc

Zonal Regional Multi-Regional Global



- Batch MapReduce processing via configurable, managed Spark & Hadoop clusters
- Handles being told to scale (adding or removing nodes) even while running jobs
- Integrated with Cloud Storage, BigQuery, Bigtable, and some Stackdriver services
- “Image versioning” switches between versions of Spark, Hadoop, & other tools
- Pay directly for underlying GCE servers used in the cluster—optionally preemptible
- Pay a Cloud Dataproc management fee per vCPU-hour in the cluster
- Best for moving *existing* Spark/Hadoop setups to GCP
 - Prefer Cloud Dataflow for new data processing pipelines — “Go with the flow”

Data Flows - Processing: Big Data and IoT



Cloud
Dataflow

Zonal Regional Multi-Regional Global



- Smartly-autoscaled & fully-managed batch or stream MapReduce-like processing
- Released as open-source Apache Beam



Amazon EMR



Spark beam



Data Flows - Processing: Big Data and IoT



Cloud
Dataflow

Zonal Regional Multi-Regional Global



- Smartly-autoscaled & fully-managed batch or stream MapReduce-like processing
- Released as open-source Apache Beam
- Autoscales & dynamically redistributes lagging work, mid-job, to optimize run time
- Integrated with Cloud Pub/Sub, Datastore, BQ, Bigtable, Cloud ML, Stackdriver, etc.
- Dataflow Shuffle service for batch offloads Shuffle ops from workers for big gains
- Effectively pay for underlying worker GCE via consolidated charges
 - Pay per second for vCPUs, RAM GBs PD/PD-SSD (more for streaming)
 - Dataflow Shuffle charged for time per GB used



Data Flows - Processing: Big Data and IoT



Cloud
Datalab

Zonal Regional Multi-Regional Global



- Interactive tool for data exploration, analysis, visualization and machine learning
- Uses Jupyter Notebook
 - “[A]n open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text. Uses include: data cleaning and transformation, numerical simulation, statistical modeling, data visualization, machine learning, and much more.”
- Supports iterative development of data analysis algorithms in Python/SQL/~JS
- Pay for GCE/GAE instance hosting and storing (on PD) your notebooks
- Pay for any other resources accessed (e.g. BigQuery)



Data Flows - Processing: Big Data and IoT



Cloud
Data Studio

Zonal Regional Multi-Regional Global



- **Big Data Visualization tool for dashboards and reporting**
- **Meaningful data stories/presentations enable better business decision making**
- **Data sources include BigQuery, Cloud SQL, other MySQL, Google Sheets, Google Analytics, Analytics 360, AdWords, DoubleClick, & YouTube channels**
- **Visualizations include time series, bar charts, pie charts, tables, heat maps, geo maps, scorecards, scatter charts, bullet charts, & area charts**
- **Templates for quick start; customization options for impactful finish**
- **Familiar G Suite sharing and real-time collaboration**
- **Pay only for services accessed**

Data Flows - Processing: Big Data and IoT



Cloud
Genomics

Zonal Regional Multi-Regional Global



- Store and process genomes and related experiments
- Query complete genomic information of large research projects in seconds
- Process many genomes and experiments in parallel
- Open industry standards (e.g. from Global Alliance for Genomics and Health)
- Supports “Requester Pays” sharing

Data Flows - Processing: Dataproc or Dataflow

Processing large-scale data

Large-scale data processing typically involves reading data from source systems such as Cloud Storage, Cloud Bigtable, or Cloud SQL, and then conducting complex normalizations or aggregations of that data. In many cases, the data is too large to fit on a single machine so frameworks are used to manage distributed compute clusters and to provide software tools that aid processing.



Cloud Dataproc



Cloud Dataflow



Cloud Dataprep

- | | | |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------|
| <ul style="list-style-type: none">• Existing Hadoop/Spark Applications• Machine Learning / Data Science Ecosystem• Tunable Cluster Parameters | <ul style="list-style-type: none">• New Data Processing Pipelines• Unified Streaming & Batch• Fully-Managed, No-Ops | <ul style="list-style-type: none">• UI-Driven Data Preparation• Scales On-Demand• Fully-Managed, No-Ops |
|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------|



Data Flows - Processing: Dataproc or Dataflow

Recommended Workloads

WORKLOADS	CLOUD DATAPROC	CLOUD DATAFLOW
Stream processing (ETL)		✓
Batch processing (ETL)	✓	✓
Iterative processing and notebooks	✓	
Machine learning with Spark ML	✓	
Preprocessing for machine learning		✓ (with Cloud ML Engine)

Hands-on lab

Data Flows - Storage: Local SSD



Local
SSD

Zonal Regional Multi-Regional Global



- Very fast 375GB solid state drives physically attached to the server



EC2
Instance
Store Vols.



Direct-Attached
Storage (DAS)

Data Flows - Storage: Local SSD



Local
SSD

Zonal Regional Multi-Regional Global



- Very fast 375GB solid state drives physically attached to the server
- Can stripe across eight of them (3TB) for even better performance
- DATA WILL BE LOST whenever the instance shuts down
 - But can survive a Live Migration
- Like all data at rest, always encrypted

Data Flows - Storage: Persistent Disk (PD)



Persistent
Disk

Zonal Regional Multi-Regional Global



- Flexible, *block-based* network-attached storage; boot disk for every GCE instance
- Perf scales with volume size; max way below Local SSD, but still plenty fast
- Persistent disks *persist*, and are replicated (zone or region) for durability
- Can resize while in use (up to 64TB), but will need file system update within VM
- Snapshots (and machine images) add even more capability and flexibility
 - “Magical”: Pay for incremental (\$ and time), but use/delete like full backups
- Not file-based NAS, but can mount to multiple instances if *all* are read-only
- Pay for GB/mo provisioned depending on perf. class; plus snapshot GB/mo used

Data Flows - Storage: Cloud Filestore



Cloud
Filestore

Zonal Regional Multi-Regional Global



- Fully-managed file-based storage



**Elastic File
System (EFS)**



**Network-Attached
Storage (NAS)**

Data Flows - Storage: Cloud Filestore



Cloud
Filestore

Zonal Regional Multi-Regional Global



- Fully-managed file-based storage
- “Predictably fast performance for your file-based workloads”
- Accessible to GCE and GKE through your VPC, via NFSv3 protocol
- Primary use case is application migration to cloud (“lift and shift”)
- Fully manages file *serving*, but *not* backups
- Pay for provisioned TBs in “Standard” (slow) or “Premium” (fast) mode
- Minimum provisioned capacity of 1TB (Standard) or 2.5TB (Premium)

Data Flows - Storage: Cloud Storage



Cloud
Storage

Zonal Regional Multi-Regional Global



- Infinitely scalable, fully-managed, versioned, and highly-durable object storage
- Designed for 99.99999999% ("eleven nines") durability
- Strongly consistent (even for overwrite PUTs and DELETEs)
- Integrated site hosting and CDN functionality
- Lifecycle transitions across classes: Multi-Regional, Regional, Nearline, Coldline
 - Diffs in cost & availability (99.95%, 99.9%, 99%, 99%), not latency (no thaw delay)
- All classes have same API, so can use gsutil and gcsfuse (but beware)
- Pay for data operations & GB-months stored by class
- Nearline/Coldline: Also pay for GBs retrieved—plus early deletion fee, if < 30/90 days



Choosing your storage: Flowchart

Do you get to choose?
Value flexibility?
Open access?
Simplicity?
Pay by usage?
Scale like mad?



Yes!



Cloud
Storage

No

NAS-based
Lift-and-Shift?
Multiple Writers?

Yes

Need absolute
max performance?

No

Yes



Local
SSD



Persistent
Disk



Cloud
Filestore

Choosing a cloud storage class: Flowchart



Databases: CloudSQL



Cloud
SQL

Zonal Regional Multi-Regional Global



- Fully-managed and reliable MySQL and PostgreSQL databases



Amazon RDS



Self-Managed
MySQL

Databases: CloudSQL



Cloud
SQL

Zonal Regional Multi-Regional Global



- Fully-managed and reliable MySQL and PostgreSQL databases
- Supports automatic replication, backup, failover, etc.
- Scaling is manual (both vertically and horizontally)
- Effectively pay for underlying GCE instances and PDs
 - Plus some baked-in service fees

Databases: Cloud Spanner



Cloud
Spanner

Zonal Regional Multi-Regional Global



- “The first horizontally scalable, strongly consistent, relational database service”
 - “From 1 to hundreds or thousands of nodes”
 - “A minimum of 3 nodes is recommended for production environments.”
- Chooses Consistency and Partition-Tolerance (CP of CAP theorem)
- But still *high Availability*: SLA has 99.999% SLO (five nines) for multi-region
 - Nothing is actually 100%, really
 - Not based on fail-over
- Pay for provisioned node time (by region/multi-region) plus used storage-time



Databases: Cloud Spanner

Cloud Spanner

High Capacity Global Spanner  

Spanner nodes: 4,228

Storage: 100 TB per month

\$1,003,324,320.00

Total Estimated Cost: \$1,003,324,320.00 per 3 years

Adjust Estimate Timeframe

1 day 1 week 1 month 1 quarter 1 year 3 years

EMAIL ESTIMATE SAVE ESTIMATE



Databases: BigQuery



BigQuery

Zonal Regional Multi-Regional Global



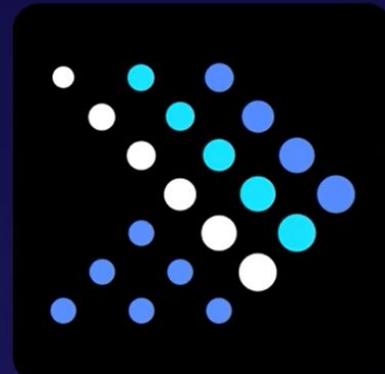
- Serverless column-store data warehouse for analytics using SQL
- Scales internally, so it “can scan TB in seconds and PB in minutes”



Amazon
Redshift



Amazon
Athena



Presto

Databases: BigQuery



BigQuery

Zonal Regional Multi-Regional Global



- Serverless column-store data warehouse for analytics using SQL
- Scales internally, so it “can scan TB in seconds and PB in minutes”
- Pay for GBs actually considered (scanned) during queries
 - Attempts to reuse cached results, which are free
- Pay for data stored (GB-months)
 - Relatively inexpensive
 - Even cheaper when table not modified for 90 days (reading still fine)
- Pay for GBs added via streaming inserts



Databases: Cloud Bigtable



Cloud
Bigtable

Zonal Regional Multi-Regional Global



- Low latency & high throughput NoSQL DB for large operational & analytical apps



DynamoDB



Cassandra

APACHE
HBASE

Databases: Cloud Bigtable



Cloud
Bigtable

Zonal Regional Multi-Regional Global



- Low latency & high throughput NoSQL DB for large operational & analytical apps
- Supports open-source HBase API
- Integrates with Hadoop, Dataflow, Dataproc
- Scales seamlessly and unlimitedly
 - Storage autoscales
 - Processing nodes must be scaled manually
- Pay for processing node hours

Databases: Cloud Datastore



Cloud
Datastore

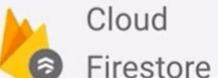
Zonal Regional Multi-Regional Global



- Managed & autoscaled NoSQL DB with indexes, queries, and ACID trans. support
- NoSQL, so queries can get complicated
 - No joins or aggregates and must line up with indexes
 - NOT, OR, and NOT EQUALS (<>, !=) operations not natively supported
- Automatic “built-in” indexes for simple filtering and sorting (ASC, DESC)
- Manual “composite” indexes for more complicated, but beware them “exploding”
- Pay for GB-months of storage used (including indexes)
- Pay for IO operations (deletes, reads, writes) performed (i.e. no pre-provisioning)



Databases: Firebase

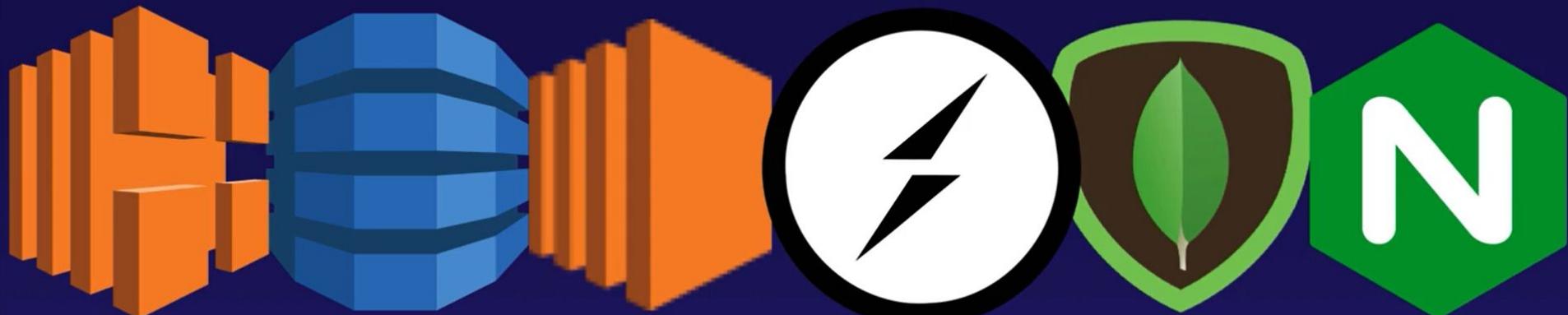


Zonal Regional Multi-Regional Global

Zonal Regional Multi-Regional Global



- NoSQL document stores with ~real-time client updates via managed websockets



ELB

DynamoDB

EC2

Socket.io

MongoDB

NGINX

Databases: Cloud Spanner



Firebase
Realtime DB



Cloud
Firestore

Zonal

Regional

Multi-Regional

Global

Zonal

Regional

Multi-Regional

Global

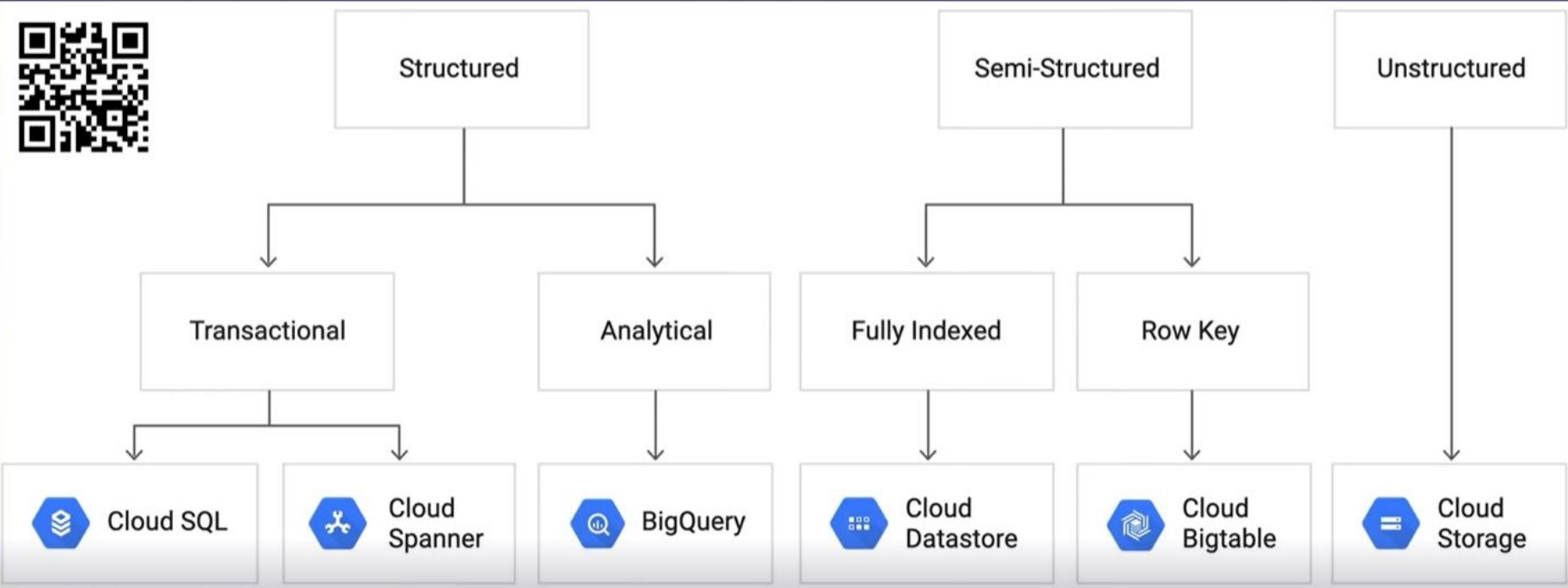


- NoSQL document stores with ~real-time client updates via managed websockets
- Firebase DB is single (potentially huge) JSON doc, located only in central US
- Cloud Firestore has collections, documents, and contained data
- Free tier (Spark), flat tier (Flame), or usage-based pricing (Blaze)
 - Realtime DB: Pay more for GB/month stored and GB downloaded
 - Firestore: Pay for operations and much less for storage and transfer

Choosing a database: Flowchart



Choosing a database: Data-first approach



Hands-on lab

Network Layer Design

Data transfer

Data Transfer: Data Transfer Appliance



Data Transfer
Appliance



- Rackable, high-capacity storage server to physically ship data to GCS



AWS Snowball

Data Transfer: Data Transfer Appliance



Data Transfer
Appliance



- **Rackable, high-capacity storage server to physically ship data to GCS**
- **Ingest only; not a way to avoid egress charges**
- **100TB or 480TB versions**
- **480TB/week is faster than a saturated 6Gbps link**

Data Transfer: Storage Transfer Service



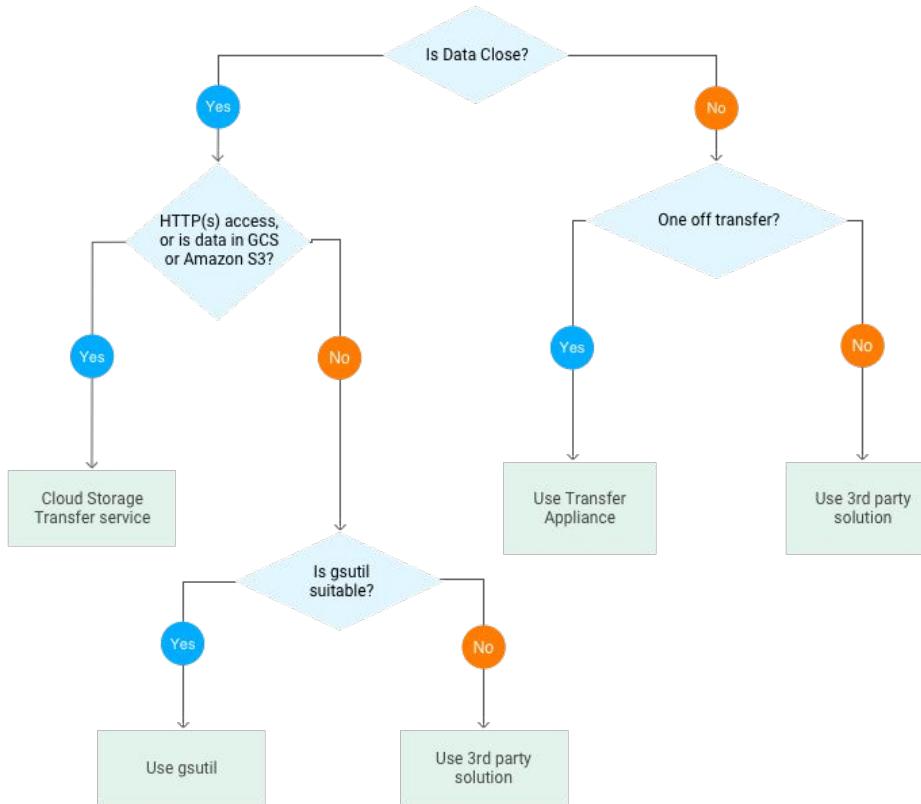
Storage Transfer
Service

Zonal Regional Multi-Regional Global



- Copies objects for you, so you don't need to set up a machine to do it
- Destination is always GCS bucket
- Source can be S3, HTTP/HTTPS endpoint, or another GCS bucket
- One-time or scheduled recurring transfers
- Free to use, but you pay for its actions

Data transfer options : Flowchart



Data Transfer: Google Domains



Google
Domains

Zonal Regional Multi-Regional Global



- **Google's registrar for domain names**
- **Private Whois records**
- **Built-in DNS or custom nameservers**
- **Supports DNSSEC**
- **Email forwarding with automatic setup of SPF and DKIM (for built-in DNS)**

Data Transfer: Cloud DNS



Cloud
DNS

Zonal Regional Multi-Regional Global



- Scalable, reliable, & managed authoritative Domain Name System (DNS) service



Amazon
Route 53



Dyn

Data Transfer: Cloud DNS



Cloud
DNS

Zonal Regional Multi-Regional Global



- Scalable, reliable, & managed authoritative Domain Name System (DNS) service
- 100% uptime guarantee
- Public and private managed zones
- Low latency globally
- Supports DNSSEC
- Manage via UI, CLI, or API
- Pay fixed fee per managed zone to store and distribute DNS records
- Pay for DNS lookups (i.e. usage)

Data Transfer: Static IP



Static IP

Zonal Regional Multi-Regional Global



- Reserve static IP addresses in projects and assign them to resources
- Regional IPs used for GCE instances & Network Load Balancers
- Global IPs used for global load balancers:
 - HTTP(S), SSL proxy, and TCP proxy
 - “Anycast IP” simplifies DNS
- Pay for reserved IPs that are not in use, to discourage wasting them

Data Transfer: Load Balancing



Load
Balancing

Zonal Regional Multi-Regional Global



- High-perf, scalable traffic distribution integrated with autoscaling & Cloud CDN



Elastic Load
Balancing



NGINX

Data Transfer: Load Balancing



Load
Balancing

Zonal

Regional

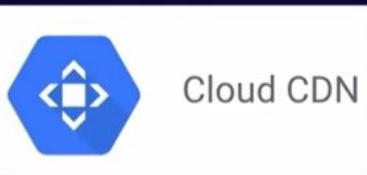
Multi-Regional

Global



- High-perf, scalable traffic distribution integrated with autoscaling & Cloud CDN
- SDN naturally handles spikes without any prewarming; no instances or devices
- Regional Network Load Balancer: health checks, round robin, session affinity
 - Forwarding rules based on IP, protocol (e.g. TCP, UDP), and (optionally) port
- Global load balancers w/ multi-region failover for HTTP(S), SSL proxy, & TCP proxy
 - Prioritize low-latency connection to region near user, then gently fail over in bits
 - Reacts quickly (unlike DNS) to changes in users, traffic, network, health, etc.
- Pay by making ingress traffic billable (cheaper than egress) plus hourly per rule

Data Transfer: Cloud CDN



Cloud CDN

Zonal Regional Multi-Regional Global



- Low-latency content delivery based on HTTP(S) CLB & integrated w/ GCE & GCS



Amazon
CloudFront



Cloudflare

Data Transfer: Cloud CDN



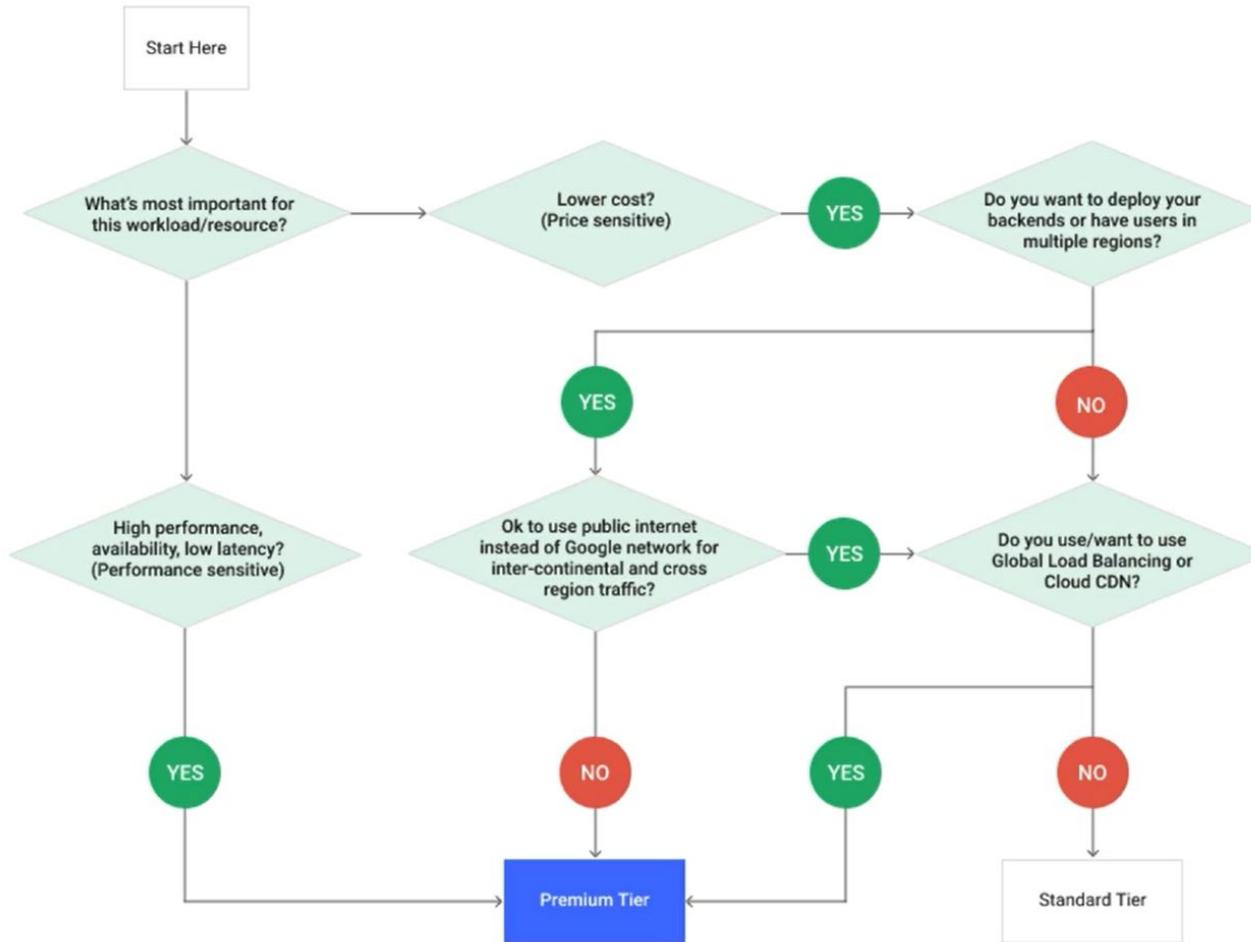
Cloud CDN

Zonal Regional Multi-Regional Global

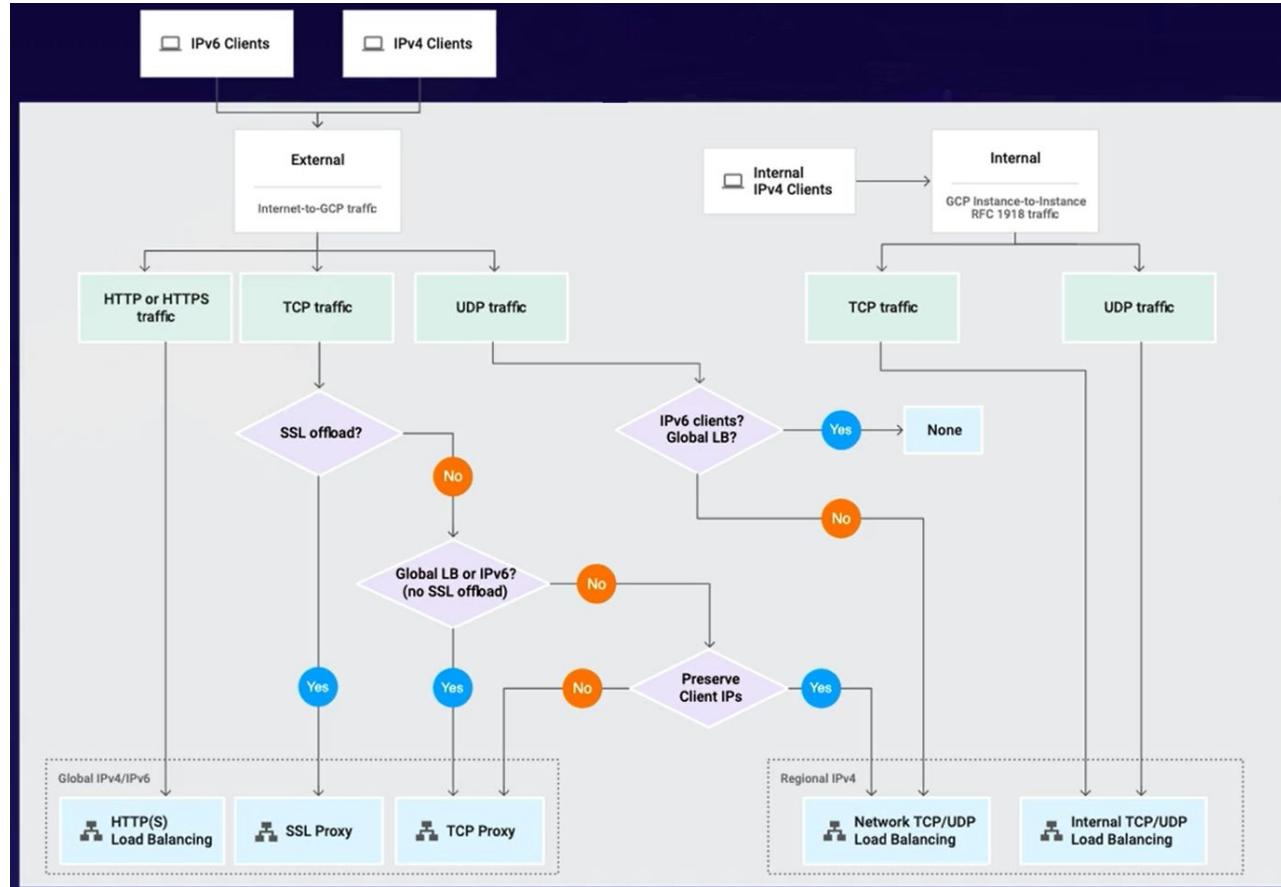


- Low-latency content delivery based on HTTP(S) CLB & integrated w/ GCE & GCS
- Supports HTTP/2 and HTTPS, but no custom origins (GCP only)
- Simple checkbox on HTTP(S) Load Balancer config turns this on
- On cache miss, pay origin→POP “cache fill” egress charges (cheaper for in-region)
- Always pay POP→client egress charges, depending on location
- Pay for HTTP(S) request volume
- Pay per cache invalidation request (not per resource invalidated)
- Origin costs (e.g. CLB, GCS) can be much lower because cache hits reduce load

Network Tier Choice : Flowchart



Load Balancer Options : Flowchart



Internal Networking: VPC



VPC

Zonal Regional Multi-Regional Global



- Global IPv4 unicast Software-Defined Network (SDN) for GCP resources
- Automatic mode is easy; custom mode gives control
- Configure subnets (each with a private IP range), routes, firewalls, VPNs, BGP, etc.
- VPC is global and subnets are regional (not zonal!)
- Can be shared across multiple projects in same org and peered with other VPCs
- Can enable private (internal IP) access to some GCP services (e.g. BQ, GCS)
- Free to configure VPC (container)
- Pay to use certain services (e.g. VPN) and for network egress

Internal Networking: Cloud Interconnect



Cloud
Interconnect

Zonal Regional Multi-Regional Global



- Options for connecting external networks to Google's network
- Private connections to VPC via Cloud VPN or Dedicated/Partner Interconnect
- Public Google services (incl. GCP) accessible via External Peering (no SLAs)
 - Direct Peering for high volume
 - Carrier Peering via a partner for lower volume
- Significantly lower egress fees

Internal Networking: Cloud VPN



Cloud
VPN

Zonal Regional Multi-Regional Global



- IPsec VPN to connect to VPC via public internet for low-volume data connections



AWS VPN



Internal Networking: Cloud VPN



Cloud
VPN

Zonal Regional Multi-Regional Global



- IPsec VPN to connect to VPC via public internet for low-volume data connections
- For persistent, static connections between gateways (i.e. not for a dynamic client)
 - Peer VPN gateway must have static (unchanging) IP
- Encrypted link to VPC (as opposed to Dedicated Interconnect), into one subnet
- Supports both static and dynamic routing
- 99.9% availability SLA
- Pay per tunnel-hour

Internal Networking: Dedicated Interconnect



Dedicated
Interconnect

Zonal Regional Multi-Regional Global



- Direct physical link between VPC and on-prem for high-volume data connections
- VLAN attachment is private connection to VPC in one region; no public GCP APIs
 - Region chosen from those supported by particular Interconnect Location
- Links are private but not encrypted; can layer your own encryption
- Redundant connections in different locations recommended for critical apps
 - Redundancy achieves 99.99% availability; otherwise 99.9% SLA
- Pay fee per 10 Gbps link, plus (relatively small) fee per VLAN attachment
- Pay reduced egress rates from VPC through Dedicated Interconnect

Internal Networking: Cloud Router



Cloud
Router

Zonal Regional Multi-Regional Global



- Dynamic routing (BGP) for hybrid networks linking GCP VPCs to external networks
- Works with Cloud VPN and Dedicated Interconnect
- Automatically learns subnets in VPC and announces them to on-prem network
- Without Cloud Router you must manage static routes for VPN
 - Changing the IP addresses on either side of VPN requires recreating it
- Free to set up
- Pay for usual VPC egress

Internal Networking: CDN Interconnect



CDN
Interconnect

Zonal Regional Multi-Regional Global



- Direct, low-latency connectivity to certain CDN providers, with cheaper egress
- For external CDNs, not Google's Cloud CDN service
 - Supports Akamai, Cloudflare, Fastly, and more
- Works for both pull and push cache fills
 - Because it's for all traffic with that CDN
- Contact CDN provider to set up for GCP project and which regions
- Free to enable, then pay less for the egress you configured

Hands-on lab

Design for security

Building and Managing

Identity and Access - Core Security



Roles

Zonal Regional Multi-Regional Global



- Roles are collections of Permissions to use or manage GCP resources



AWS IAM Policies

Identity and Access - Roles



Roles

Zonal Regional Multi-Regional Global



- Roles are **collections of Permissions** to use or manage GCP resources
- Permissions allow you to perform certain actions: Service . Resource . Verb
- Primitive Roles: Owner, Editor, Viewer
 - Viewer is read-only; Editor can change things; Owner can control access & billing
 - Pre-date IAM service, may still be useful (e.g. dev/test envs), but often too broad
- Predefined Roles: Give granular access to specific GCP resources (IAM)
 - E.g.: roles/bigquery.dataEditor, roles/pubsub.subscriber
- Custom Roles: Project- or Org-level collections you define of granular permissions

Identity and Access - Cloud IAM



Cloud IAM

Zonal Regional Multi-Regional Global



- Controls access to GCP resources: authorization, not really authentication/identity



AWS IAM

Identity and Access - Cloud IAM



Cloud IAM

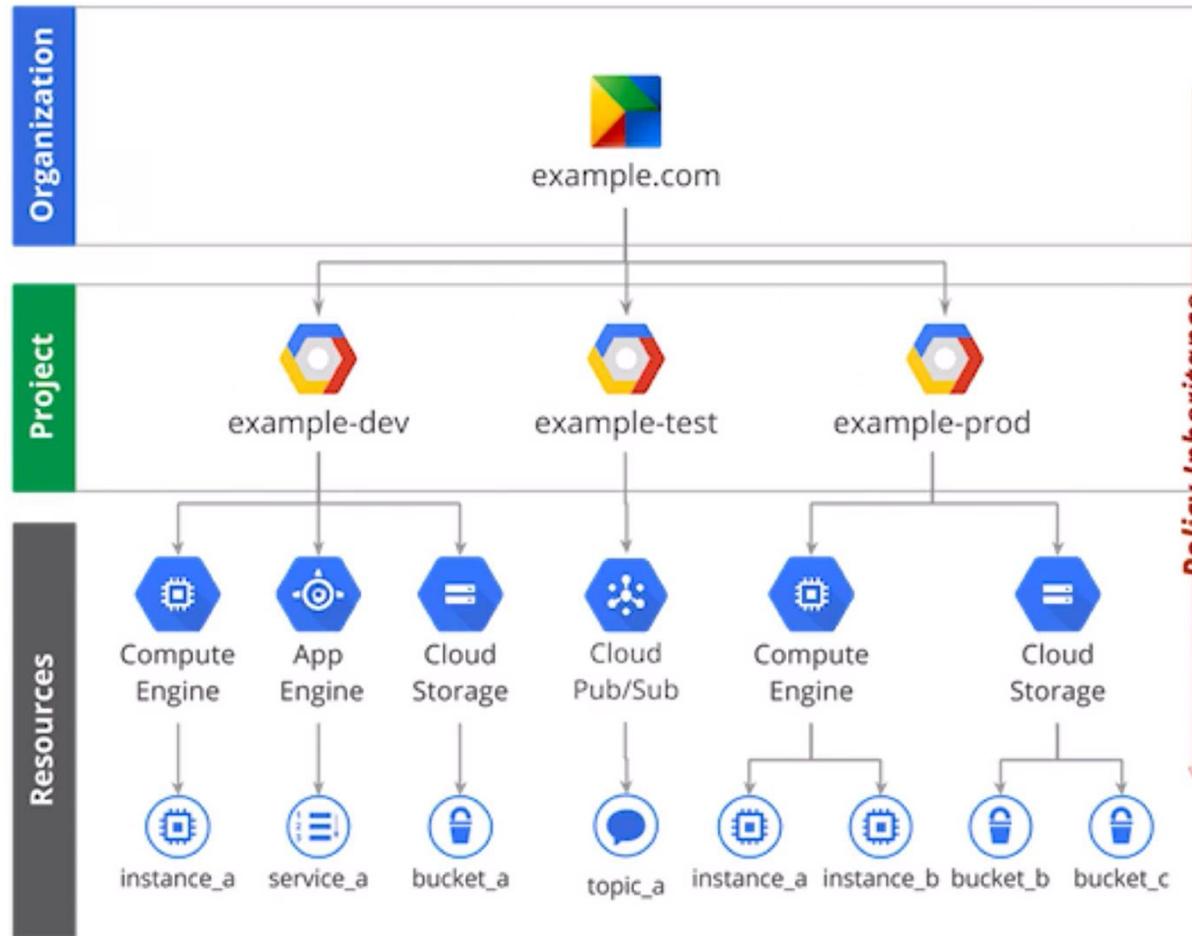
Zonal Regional Multi-Regional Global



- Controls access to GCP resources: authorization, not really authentication/identity
- Member is user, group, domain, service account, or the public (e.g. "allUsers")
 - Individual Google account, Google group (👍), G Suite / Cloud Identity domain
 - Service account (👍) belongs to application/instance, not individual end user
 - Every identity has a unique e-mail address, including service accounts
- Policies bind Members to Roles at a hierarchy level: Org, Folder, Project, Resource
 - Answer: *Who can do what to which thing(s)?*
- IAM is free; pay for authorized GCP service usage



Identity and Access - Core Security



Identity and Access - Service Accounts



Service
Accounts

Zonal Regional Multi-Regional Global



- Special type of Google account that represents an application, not an end user
- Can be “assumed” by applications or individual users (when so authorized)
- “Important: For almost all cases, whether you are developing locally or in a production application, you should use service accounts, rather than user accounts or API keys.”
- Consider resources and permissions required by application; use least privilege
- Can generate and download private keys (user-managed keys), for non-GCP, but...
- Cloud-Platform-managed keys (👍) are better, for GCP (i.e. GCF, GAE, GCE, and GKE)
 - No direct downloading: Google manages private keys & rotates them once a day

Identity and Access - Cloud Identity



Cloud
Identity

Zonal Regional Multi-Regional Global



- Identity as a Service (IDaaS, not DaaS) to provision and manage users and groups



AWS IAM



G Suite



GMail/Google Account



Active Directory

Identity and Access - Core Security



Cloud
Identity

Zonal Regional Multi-Regional Global



- Identity as a Service (**IDaaS**, not DaaS) to provision and manage users and groups
- Free Google Accounts for non-G-Suite users, tied to a verified domain
- Centrally manage all users in Google Admin console; supports compliance
- 2-Step verification (2SV/MFA) and enforcement (👍), including security keys
- Sync from Active Directory and LDAP directories via Google Cloud Directory Sync
- Identities work with other Google services (e.g. Chrome)
- Identities can be used to SSO with other apps via OIDC, SAML, OAuth2
- Cloud Identity is free; pay for authorized GCP service usage

Identity and Access - Security Key Enforcement



Security Key
Enforcement

Zonal Regional Multi-Regional Global



- **USB or Bluetooth 2-step verification device that prevents phishing**
- **Not like just getting a code via email or text message...**
- **Device also verifies the target service**
- **Eliminates man-in-the-middle (MITM) attacks against GCP credentials**

Identity and Access - Resource Manager



Resource
Manager

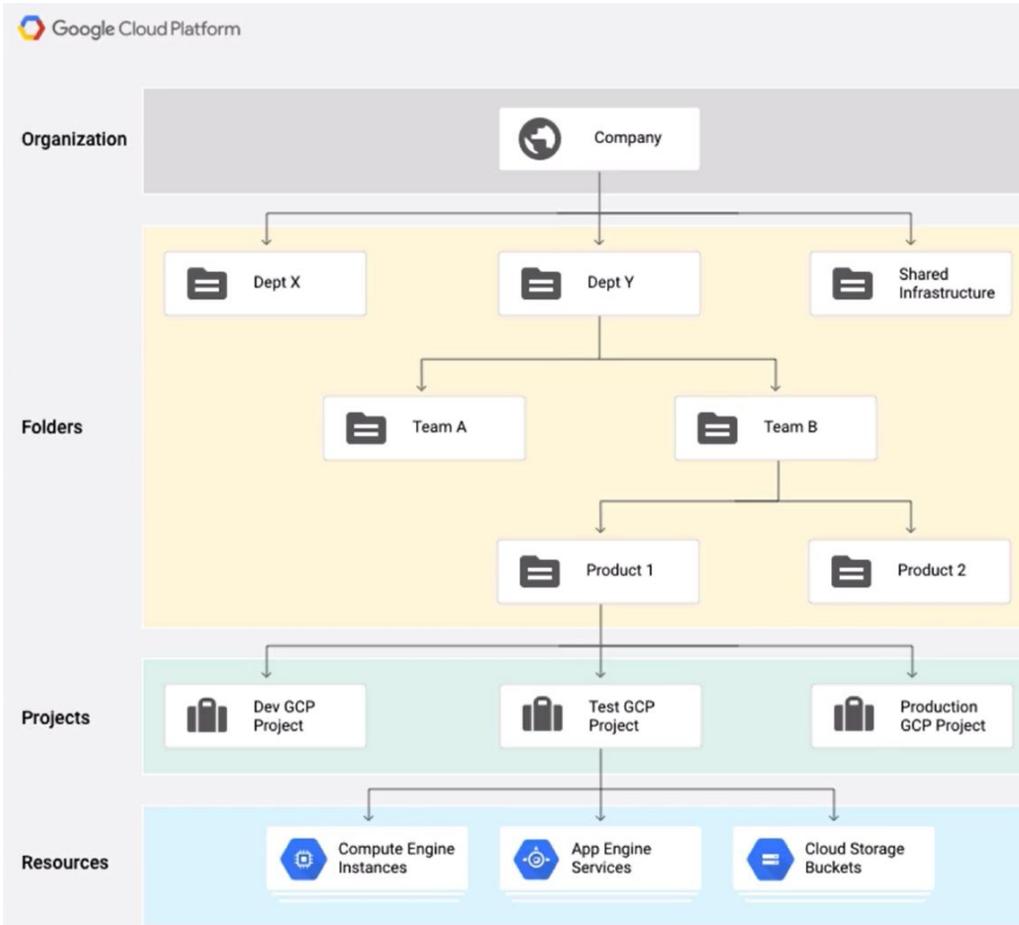
Zonal Regional Multi-Regional Global



- Centrally manage & secure organization's projects with custom folder hierarchy
- Organization resource is root node in hierarchy; folders per your business needs
- Tied 1:1 to a Cloud Identity / G Suite domain, then owns all newly-created projects
 - Without this organization, specific identities (people) must own GCP projects
- "Recycle bin" allows undeleting projects
- Define custom IAM roles at org level
- Apply IAM policies at organization, folder, or project levels



Identity and Access - Resource Manager



Identity and Access - Cloud IAP

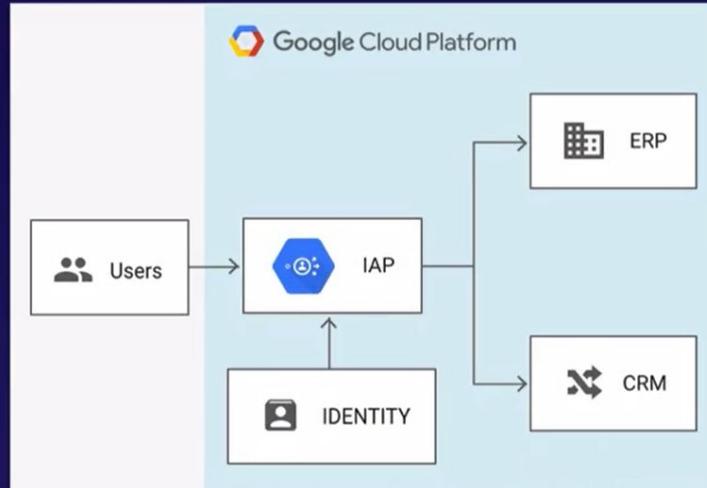


Cloud IAP

Zonal Regional Multi-Regional Global



- Guards apps running on GCP via identity verification, not VPN access
- Based on CLB & IAM, and only passes authed requests through



Identity and Access - Cloud IAP



Cloud IAP

Zonal Regional Multi-Regional Global



- Guards apps running on GCP via identity verification, not VPN access
- Based on CLB & IAM, and only passes authed requests through



Amazon API
Gateway

Identity and Access - Cloud Audit Logging



Cloud Audit
Logging

Zonal Regional Multi-Regional Global



- Answers the questions "Who did what, where, and when?" within GCP projects
- Maintains non-tamperable audit logs for each project and organization:
 - Admin Activity and System Events (400 day retention)
 - Access Transparency (400 day retention)
 - Shows actions by Google support staff
 - Data Access (30 day retention)
 - For GCP-visible services (e.g. Can't see into MySQL DB on GCE)
 - Data Access logs priced through Stackdriver Logging; rest are free

Security Management : Cloud Armor



Cloud
Armor

Zonal Regional Multi-Regional Global



- Edge-level protection from DDoS & other attacks on global HTTP(S) LB
- Offload work: Blocked attacks never reach your systems
- Monitor: Detailed request-level logs available in Stackdriver Logging
- Manage IPs with CIDR-based allow/block lists (aka whitelist/blacklist)
- More intelligent rules forthcoming (e.g. XSS, SQLi, geo-based, custom)
- Preview effect of changes before making them live
- Pay per policy and rule configured, plus for incoming request volume

Security Management: Security Scanner



Security
Scanner

Zonal Regional Multi-Regional Global



- Free but limited GAE app vulnerability scanner with “very low false positive rates”



Amazon
Inspector



Trustwave
App Scanner



Qualys Web
Application Scanning

Security Management: Security Scanner



Security
Scanner

Zonal Regional Multi-Regional Global



- Free but limited GAE app vulnerability scanner with “very low false positive rates”
- “After you set up a scan, Cloud Security Scanner automatically crawls your application, following all links within the scope of your starting URLs, and attempts to exercise as many user inputs and event handlers as possible.”
- Can identify:
 - Cross-site-scripting (XSS)
 - Flash injection
 - Mixed content (HTTP in HTTPS)
 - Outdated/insecure libraries

Security Management: Cloud DLP API



Cloud
DLP API

Zonal Regional Multi-Regional Global



- Finds and optionally redacts sensitive info in unstructured data streams
- Helps you minimize what you collect, expose, or copy to other systems
- 50+ sensitive data detectors, including: credit card numbers, names, social security numbers, passport numbers, driver's license numbers (US and some other jurisdictions), phone numbers, and other personally identifiable information (PII)
- Data can be sent directly, or API can be pointed at GCS, BQ, or Cloud DataStore
- Can scan both text and images
- Pay for amount of data processed (per GB)—and gets cheaper when large volume
 - Pricing for storage now very simple (June 2019), but for streaming is still a mess

Security Management: Event Threat Detection



Event Threat
Detection

Zonal Regional Multi-Regional Global



- Automatically scans your Stackdriver logs for suspicious activity
- Uses industry-leading threat intelligence, including Google Safe Browsing
- Quickly detects many possible threats, including:
 - Malware, cryptomining, outgoing DDoS attacks, port scanning, brute-force SSH
 - Also: Unauthorized access to GCP resources via abusive IAM access
- Can export parsed logs to BigQuery for forensic analysis
- Integrates with SIEMs like Google's Cloud SCC or via Cloud Pub/Sub
- No charge for ETD, but charged for its usage of other GCP services (like SD Logging)

Security Management: Cloud SCC



Cloud SCC

Zonal Regional Multi-Regional Global



- “Comprehensive security management and data risk platform for GCP”
- Security Information and Event Management (SIEM) software



AWS
Security Hub



Splunk ES
(Enterprise Security)



Sumo Logic

Security Management: Cloud SCC

Security Command Center + ADD SECURITY SOURCES SETTINGS

DASHBOARD ASSET FINDINGS

Assets 1 DAY Findings

Type	Deleted	New	Total
All	2	23	500
Organization	3	3	50
Project	0	10	40
Application	0	1	30
Service	0	0	30
Address	0	0	20
Disk	0	0	10
Firewall	0	23	5
instance	2	3	4
Network	3	1	3
Route	2	3	2
Subnetwork	1	4	1
Kind	2	3	1
Bucket	3	4	1

VIEW ASSET INVENTORY

Findings Summary
631 total security findings

Source	Count	Type	Count
Event Threat Detection	374	RedLock	10
Security Health Analytics	112	Cloudflare	10
Enterprise Phishing Protection	15	Qualys	8
Crowdstrike	14	Data Loss Prevention	7
Palo Alto Networks	12	+10 more	

Event Threat Detection
374 total security findings

Active threats (last 24 hours)			Active threats (last 7 days)		
Threat	Severity	Count	Type	Severity	Count
Malware: domain		8	Malware: domain		52
Cryptomining: IP		4	Malware: IP		37
Malware: hash		4	Malware: hash		32
Brute force: SSH		2	IAM: anomalous grant		11
+4 more			+4 more		

Security Management: Cloud SCC



Cloud SCC

Zonal Regional Multi-Regional Global



- “Comprehensive security management and data risk platform for GCP”
- Security Information and Event Management (SIEM) software
- “Helps you prevent, detect, & respond to threats from a single pane of glass”
- Use “Security Marks” (aka “marks”) to group, track, and manage resources
- Integrate ETD, Cloud Scanner, DLP, & many external security finding sources
- Can alert to humans & systems; Can export data to external SIEM
- Free! But charged for services used (e.g. DLP API, if configured)
- Could also be charged for excessive uploads of external findings



Encryption Key Management: Cloud KMS



Cloud KMS

Zonal Regional Multi-Regional Global



- Low-latency service to manage and use cryptographic keys
- Supports symmetric (e.g. AES) and asymmetric (e.g. RSA, EC) algorithms



AWS KMS



Encryption Key Management: Cloud KMS



Cloud KMS

Zonal Regional Multi-Regional Global



- Low-latency service to manage and use cryptographic keys
- Supports symmetric (e.g. AES) and asymmetric (e.g. RSA, EC) algorithms
- Move secrets out of code (and the like) and into the environment, in a secure way
- Integrated with IAM & Cloud Audit Logging to authorize & track key usage
- Rotate keys used for new encryption either automatically or on demand
 - Still keeps old active key versions, to allow decrypting
- Key deletion has 24 hour delay, "to prevent accidental or malicious data loss"
- Pay for active key versions stored over time
- Pay for key use operations (i.e. encrypt/decrypt; admin operations are free)

Encryption Key Management: Cloud HSM



Cloud HSM

Zonal Regional Multi-Regional Global



- Cloud KMS keys managed by FIPS 140-2 Level 3 certified HSMs
- Device hosts encryption keys and performs cryptographic operations
- Enables you to meet compliance that mandates hardware environment



AWS CloudHSM



SafeNet HSM

Encryption Key Management: Cloud KMS



Cloud HSM

Zonal Regional Multi-Regional Global



- Cloud KMS keys managed by FIPS 140-2 Level 3 certified HSMs
- Device hosts encryption keys and performs cryptographic operations
- Enables you to meet compliance that mandates hardware environment
- Fully integrated with Cloud KMS
 - Same API, features, IAM integration
- Priced like Cloud KMS: Active key versions stored & key operations
 - But some key types more expensive: RSA, EC, Long AES

Encryption Key Management Options: Flowchart



Encryption
by default

MORE AUTOMATED

World-class encryption
activated by default on
GCP



Cloud Key Management
Service



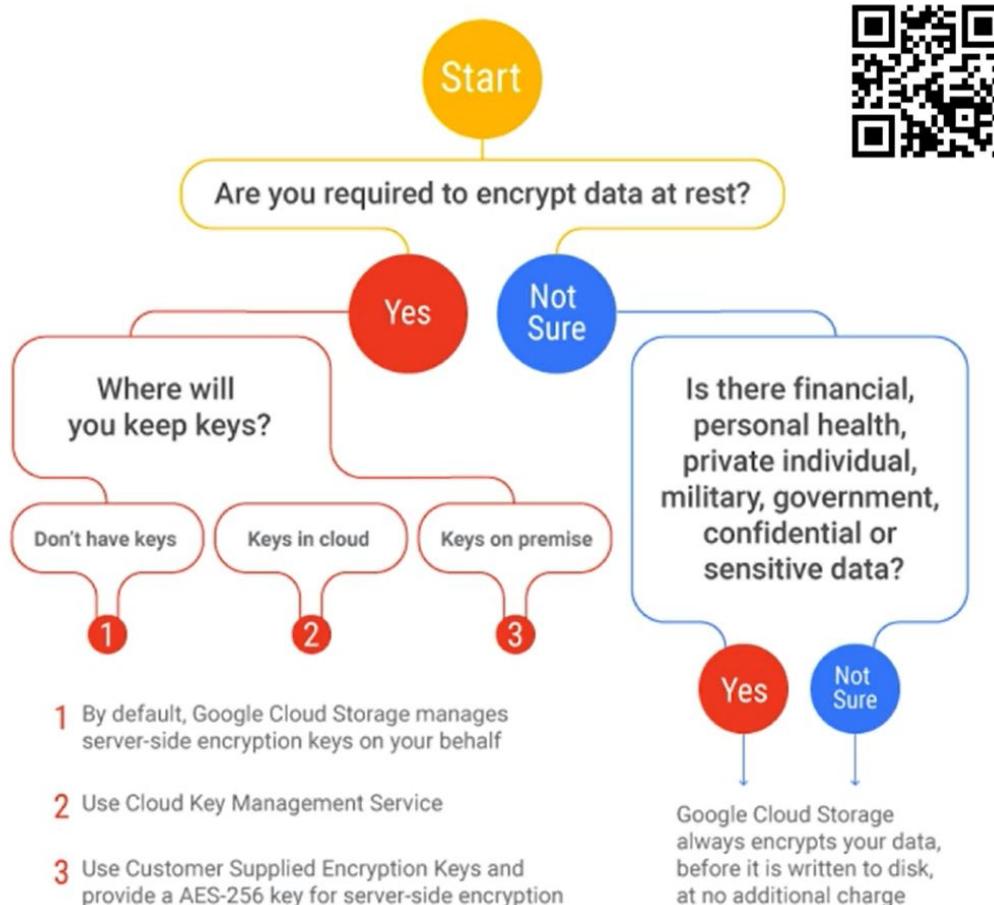
Customer-Supplied
Encryption Keys

MORE CONTROL

Keep keys in the cloud,
for direct use by cloud
services

Keep keys on-premise,
and use them to secure
your cloud services

Encryption Key Management Options: Flowchart



Operations and Management



Stackdriver
Monitoring

Zonal Regional Multi-Regional Global



- Gives visibility into perf, uptime, & overall health of cloud apps (based on collectd)
- Includes built-in/custom metrics, dashboards, global uptime monitoring, & alerts
- Follow the trail: Links from alerts to dashboards/charts to logs to traces
- Cross-cloud: GCP, of course, but monitoring agent also supports AWS
- Alerting policy config includes multi-condition rules & resource organization
- Alert via email, GCP Mobile App, SMS, Slack, PagerDuty, AWS SNS, webhook, etc.
- Automatic GCP/Anthos metrics always free
- Pay for API calls & per MB for custom or AWS metrics

Operations and Management



Stackdriver

Zonal Regional Multi-Regional Global



- Family of services for monitoring, logging, & diagnosing apps on GCP/AWS/hybrid
- Service integrations add lots of value—among Stackdriver and with GCP
- One Stackdriver account can track multiple:
 - GCP projects
 - AWS accounts
 - Other resources
- Simple usage-based pricing
 - No longer previous system of tiers, allotments, and overages

Operations and Management



Stackdriver
Monitoring

Zonal Regional Multi-Regional Global



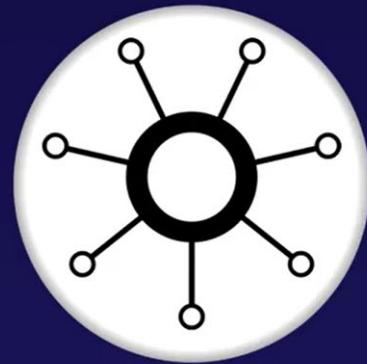
- Gives visibility into perf, uptime, & overall health of cloud apps (based on collectd)



CloudWatch Metrics
& Dashboards



DATADOG



collectd

Operations and Management



Stackdriver
Logging

Zonal Regional Multi-Regional Global



- Store, search, analyze, monitor, and alert on log data & events (based on Fluentd)
- Collection built into some GCP, AWS support with agent, or custom send via API
- Debug issues via integration with Stackdriver Monitoring, Trace & Error Reporting
- Create real-time metrics from log data, then alert or chart them on dashboards
- Send real-time log data to BigQuery for advanced analytics and SQL-like querying
- Powerful interface to browse, search, and slice log data
- Export log data to GCS to cost-effectively store log archives
- Pay per GB ingested & stored for one month, but first 50GB/project free
- Cloud Audit Logging / Access Transparency logs also free

Operations and Management



Stackdriver
Error Reporting

Zonal Regional Multi-Regional Global



- Counts, analyzes, aggregates, & tracks crashes in helpful centralized interface
- Smartly aggregates errors into meaningful groups tailored to language/framework
- Instantly alerts when a new app error cannot be grouped with existing ones
- Link directly from notifications to error details:
 - Time chart, occurrences, affected user count, first/last seen dates, cleaned stack
 - Exception stack trace parser knows Java, Python, JavaScript, Ruby, C#, PHP, & Go
 - Jump from stack frames to source to start debugging
 - No direct charge; pay for source data in Stackdriver Logging

Operations and Management



Stackdriver
Trace

Zonal Regional Multi-Regional Global



- Tracks and displays call tree & timings across distributed systems, to debug perf



AWS X-Ray



Operations and Management



Stackdriver
Trace

Zonal Regional Multi-Regional Global



- Tracks and displays call tree & timings across distributed systems, to debug perf
- Automatically captures traces from Google App Engine
- Trace API and SDKs for Java, Node.js, Ruby, and Go capture traces from anywhere
- Zipkin collector allows Zipkin tracers to submit data to Stackdriver Trace
- View aggregate app latency info or dig into individual traces to debug problems
- Generate reports on demand and get daily auto reports per traced app
- Detects app latency shift (degradation) over time by evaluating perf reports
- Pay for ingesting and retrieving trace spans

Operations and Management



Stackdriver
Debugger

Zonal Regional Multi-Regional Global



- Grabs program state (callstack, variables, expressions) in live deploys; low impact
- Logpoints repeat for up to 24h; fuller snapshots run once but can be conditional
- Source view supports Cloud Source Repository, Github, Bitbucket, local, & upload
- Java and Python supported on GCE, GKE, and GAE (Standard and Flex)
- Node.js and Ruby supported on GCE, GKE, and GAE Flex; Go only on GCE and GKE
- Automatically enabled for Google App Engine apps; agents available for others
- Share debugging sessions with others (just send URL)
- Free to use

Operations and Management



Stackdriver
Profiler

Zonal Regional Multi-Regional Global



- Continuous CPU and memory profiling to improve perf & reduce cost
- Low overhead (Typical: 0.5%; Max: 5%)—so use it in prod, too!
- Supports Go, Java, Node.js, and Python (3.2+)
- Agent-based
- Saves profiles for 30 days
- Can download profiles for longer-term storage
- Free to use

Operations and Management



Deployment
Manager

Zonal Regional Multi-Regional Global



- Create/manage resources via declarative templates: “Infrastructure as Code”
- Declarative allows automatic parallelization
- Templates written in YAML, Python, or Jinja2
- Supports input and output parameters, with JSON schema
- Create and update of deployments both support preview
- Free service; just pay for resources involved in deployments

Operations and Management



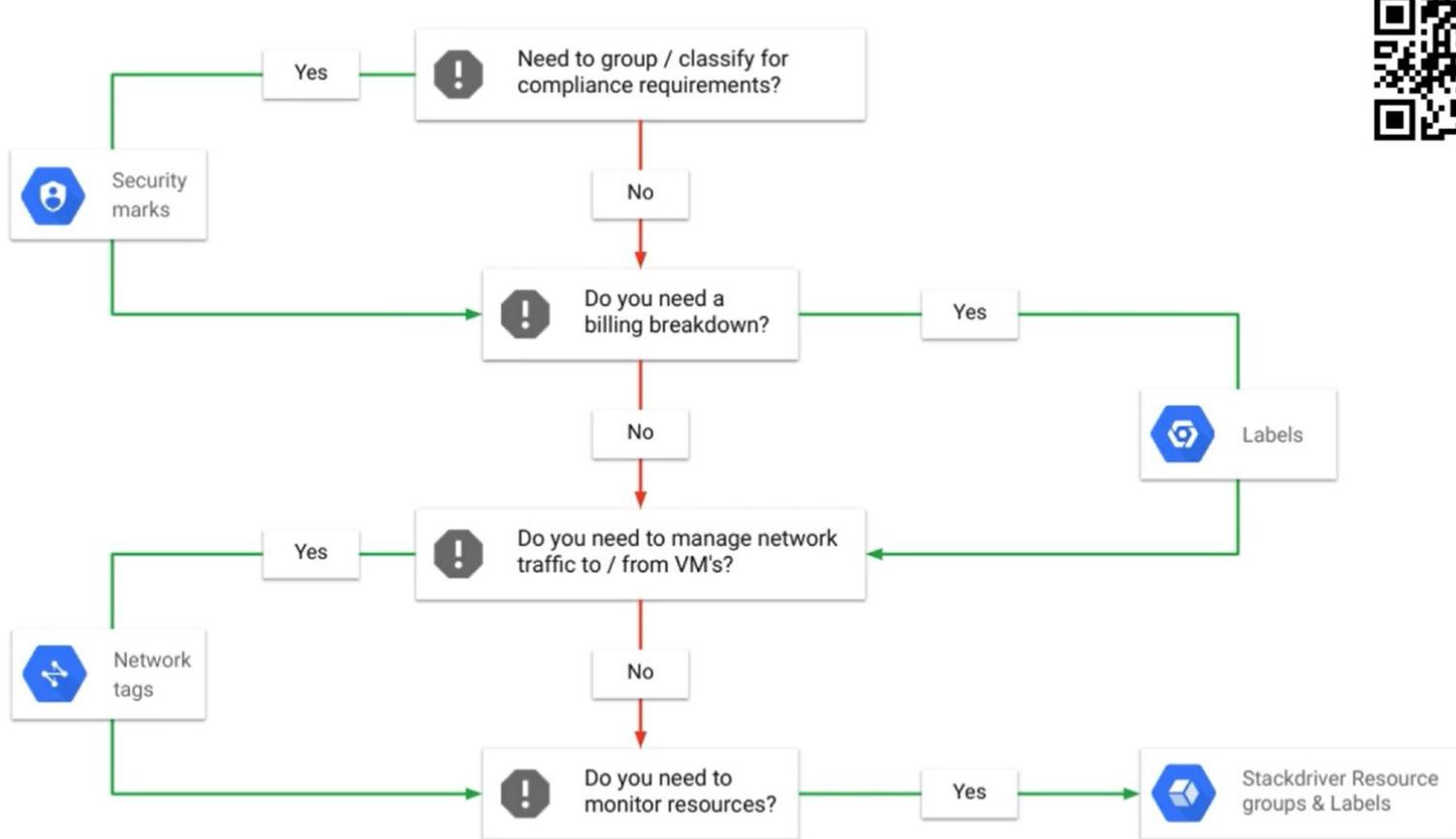
Cloud
Billing API

Zonal Regional Multi-Regional Global



- Programmatically manage billing for GCP projects and get GCP pricing
- Billing config
 - List billing accounts; get details and associated projects for each
 - Enable (associate), disable (disassociate), or change project's billing account
- Pricing
 - List billable SKUs; get public pricing (including tiers) for each
 - Get SKU metadata like regional availability
- Export of current bill to GCS or BQ is possible—but configured via console, not API

Resource Annotation Options: Flowchart



Development and APIs



Cloud Source
Repositories

Zonal Regional Multi-Regional Global



- Hosted private Git repositories, with integrations to GCP and other hosted repos
- Supports standard Git functionality
- No enhanced workflow support like pull requests
- Can set up automatic sync from GitHub or Bitbucket
- Natural integration with Stackdriver debugger for live-debugging deployed apps
- Pay per project-user active each month (not prorated)
- Pay per GB-month of data storage (prorated)
- Pay per GB of data egress

Development and APIs



Cloud Build

Zonal Regional Multi-Regional Global



- Continuously takes source code and builds, tests, and deploys it — CI/CD service
- Trigger from Cloud Source Repository (by branch, tag, or commit) or zip in GCS
 - Can trigger from GitHub and Bitbucket via Cloud Source Repositories RepoSync
- Runs many builds in parallel (currently 10 at a time)
- Dockerfile: super-simple build+push—plus scans for package vulnerabilities
- JSON/YAML file: flexible & parallel steps
- Push to GCR & export artifacts to GCS—or anywhere your build steps write
- Maintains build logs and build history
- Pay per minute of build time—but free tier is 120 minutes per day

Development and APIs



Container
Registry

Zonal Regional Multi-Regional Global



- Fast, private Docker image storage (**based on GCS**) with Docker V2 Registry API



Amazon ECR



Docker Hub

Development and APIs



Container
Registry

Zonal Regional Multi-Regional Global



- Fast, private Docker image storage (based on GCS) with Docker V2 Registry API
- Creates & manages a multi-regional GCS bucket, then translates GCR calls to GCS
- IAM integration simplifies builds and deployments within GCP
- Quick deploys because of GCP networking to GCS
- Directly compatible with standard Docker CLI; native Docker Login support
- UX integrated with Cloud Build & Stackdriver Logs
- UI to manage tags and search for images
- Pay directly for storage and egress of underlying GCS (no overhead)

Development and APIs



Cloud
Endpoints

Zonal Regional Multi-Regional Global



- Handles authorization, monitoring, logging, & API keys for APIs backed by GCP



Amazon API
Gateway



NGINX

Development and APIs



Cloud
Endpoints

Zonal Regional Multi-Regional Global



- Handles authorization, monitoring, logging, & API keys for APIs backed by GCP
- Proxy instances are distributed and hook into Cloud Load Balancer
- Super-fast Extensible Service Proxy (ESP) container based on nginx: <1 ms / call
- Uses JWTs and integrates with Firebase, Auth0, & Google Auth
- Integrates with Stackdriver Logging and Stackdriver Trace
- Extensible Service Proxy (ESP) can transcode HTTP/JSON to gRPC
 - But API needs to be resource-oriented (i.e. RESTful)
- Pay per call to your API



Development and APIs



Apigee

Zonal Regional Multi-Regional Global



- Full-featured & enterprise-scale API management platform for whole API lifecycle
- Transform calls between different protocols: SOAP, REST, XML, binary, custom
- Authenticate via OAuth/SAML/LDAP; authorize via Role-Based Access Control
- Throttle traffic with quotas, manage API versions, etc.
- Apigee Sense identifies and alerts administrators to suspicious API behaviors
- Apigee API Monetization supports various revenue models / rate plans
- Team and Business tiers are flat monthly rate with API call quotas & feature sets
- “Enterprise” tier and special feature pricing are “Contact Sales”

Hands-on lab

Thank you