

Proprietary + Confidential

**398 million**

Google

# 398 million

Google Cloud mitigated the largest DDoS attack as of 2023,  
peaking at 398 million requests per second

# What is Google Cloud?

Encryption at rest by  
default

115+ zones

Open Source

Security

Modern Infrastructure

35+ regions

Sustainable Tools and  
Infrastructure

Collaboration Cloud

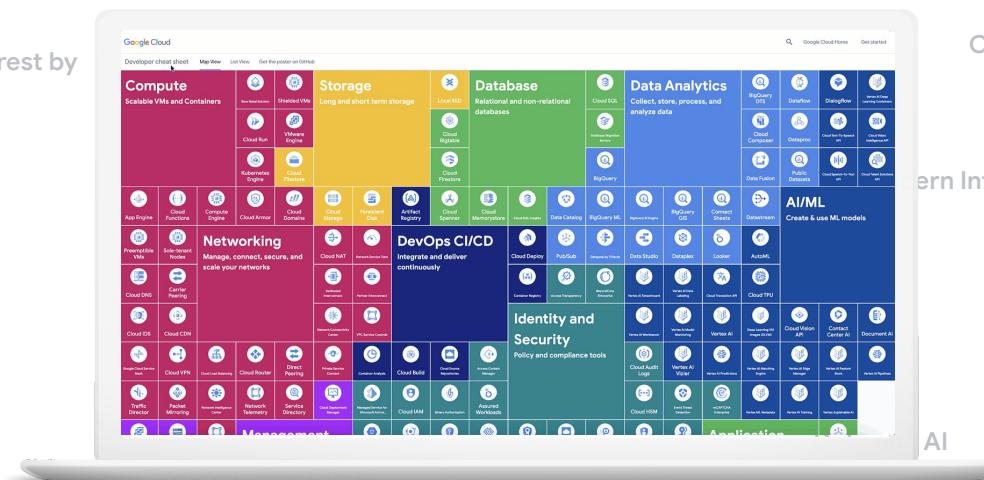
Data and AI

Google

# What is Google Cloud?

## Encryption at rest by default

35+ regions



Open Source

## Modern Infrastructure

A

Google

# What is Google Cloud?

It's how you can...

Session 1

**Host your first  
application**

Session 2

**Store and manage  
your data**

Session 3

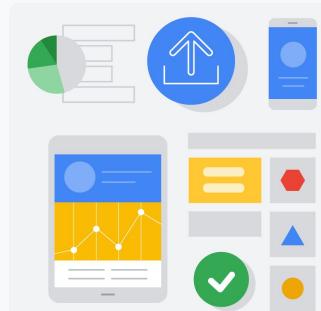
**Analyze your data**

Session 4

**Get started with  
AI**

## Click to Deploy Solutions

- Move from business use case → deployment in **one-click**
- Covers a wide range of solutions from different pillars incl.
  - Application Modernisation
  - Data Analytics/Management,
  - AI
  - Security and Monitoring
- Reflects Google Cloud best practices and model architectures
- Will be used throughout this webinar to demo our solutions and for you to get started building!



[goo.gle/loff-demo24-click-to-deploy](https://goo.gle/loff-demo24-click-to-deploy)

Google

# How to host your first application

Google Cloud

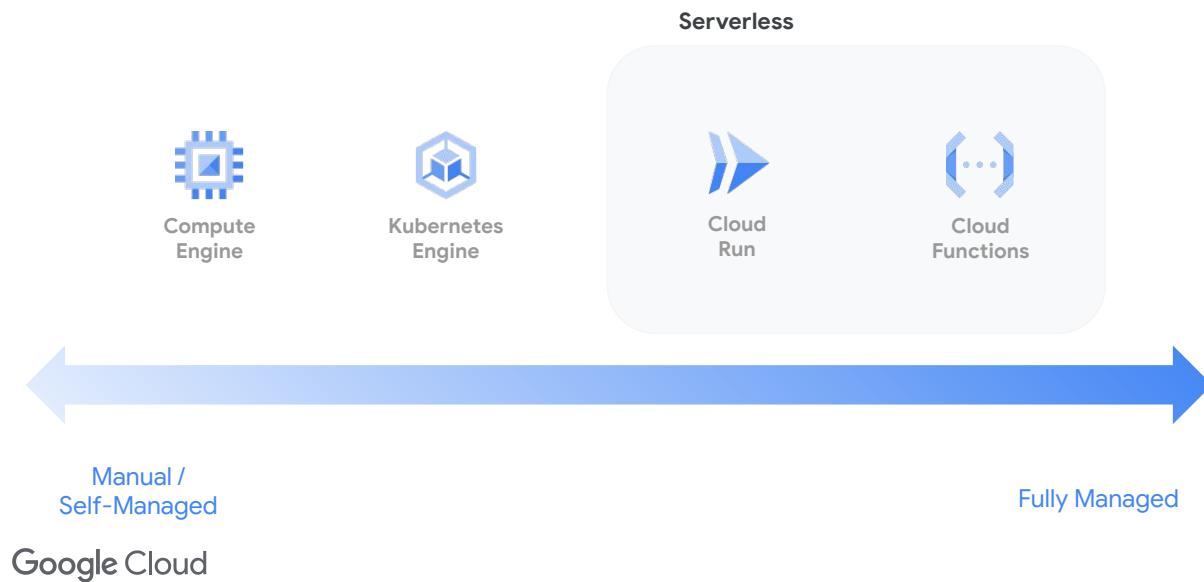


**What are you building?**

Google Cloud

# Step 1: Choose your platform

## Compute on Google Cloud

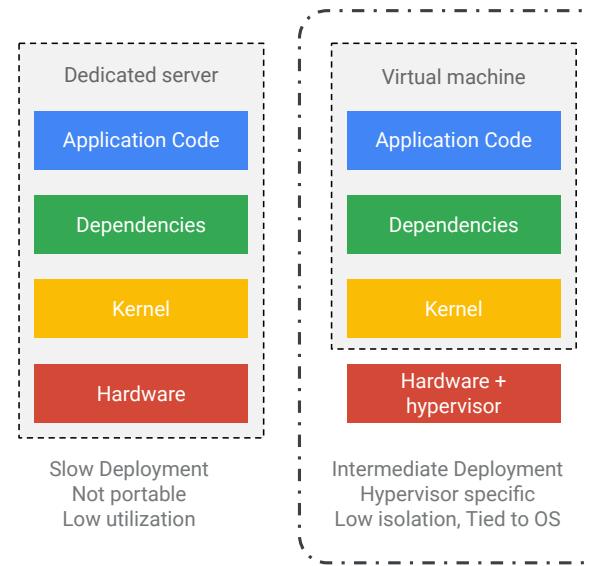
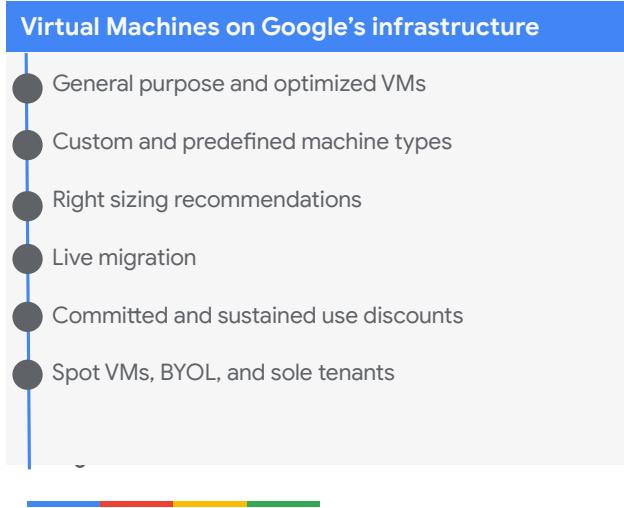


# Compute Engine

Proprietary + Confidential

## What is a Virtual Machine?

A VM is a virtualized instance of a computer that can perform almost all of the same functions as a computer, including running applications and operating systems.

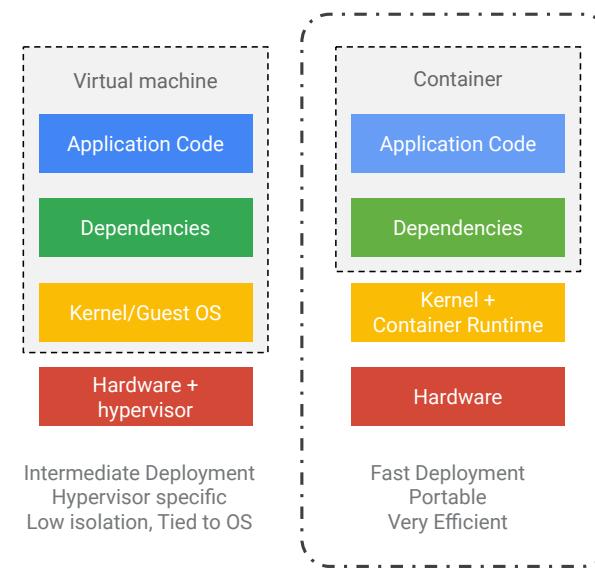
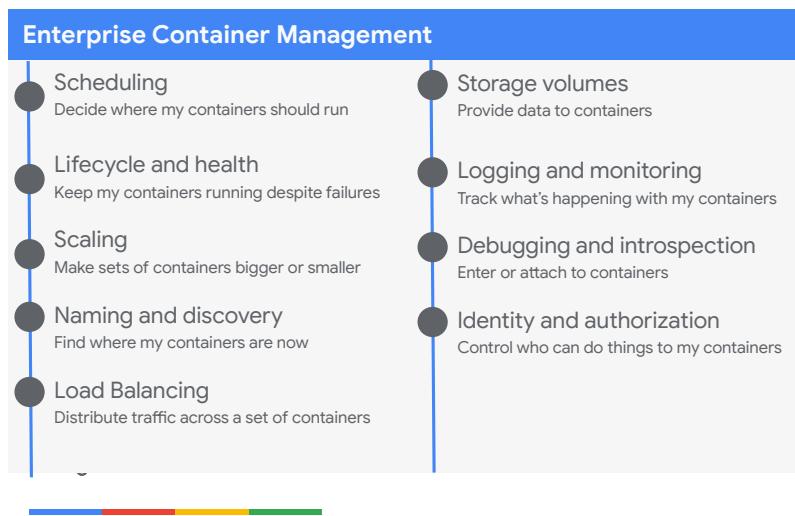


# Kubernetes Engine

Proprietary + Confidential

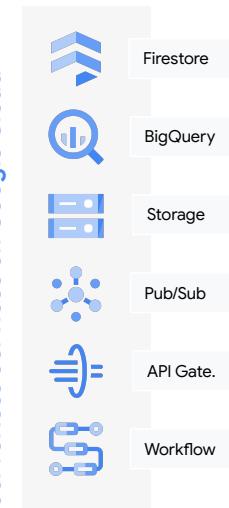
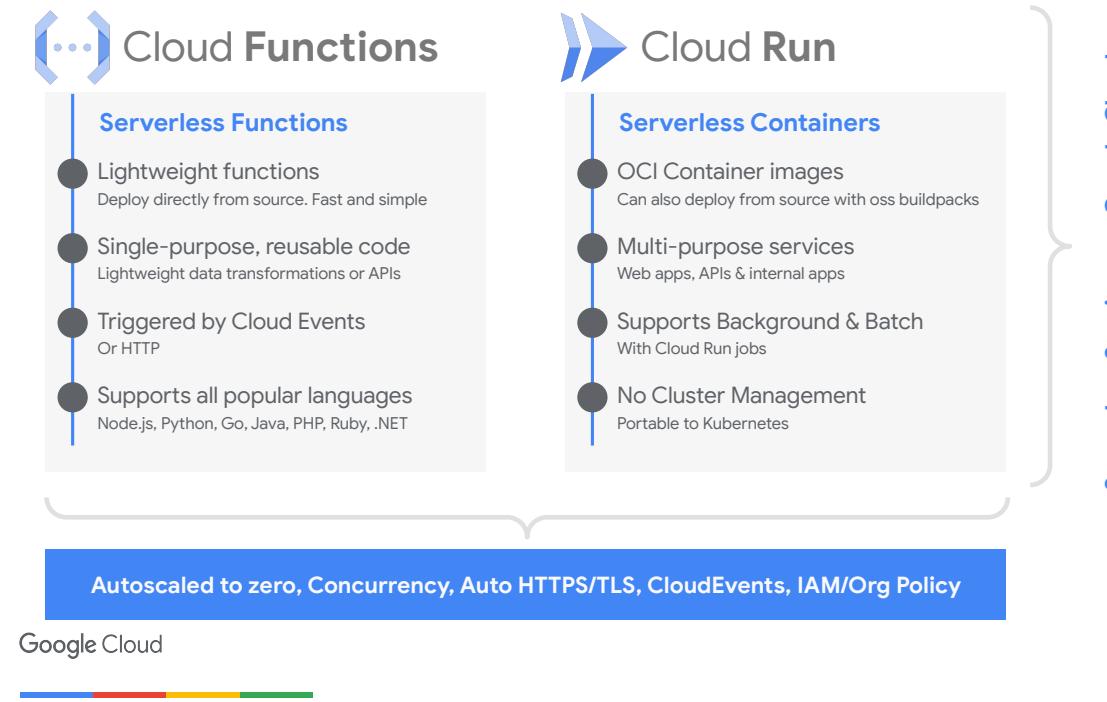
## What is a Container?

Containers are lightweight packages of your application code together with dependencies such as specific versions of programming language runtimes and libraries required to run your software services.

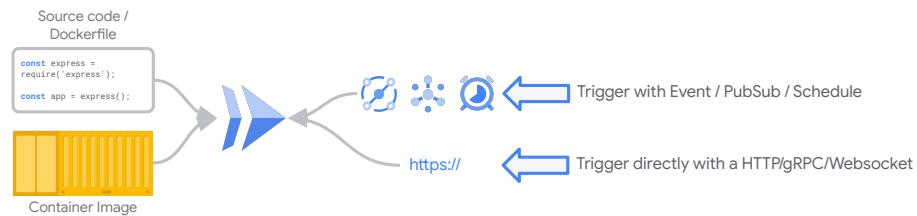


# Serverless Portfolio

Proprietary + Confidential



# Cloud Run



## Web Services, APIs, Microservices

- Optimized for developer productivity  
Fast time to market - minimal configuration and no knowledge of Kubernetes necessary
- Automatically managed & scaled  
Hyper-elastic scaling, including scale to zero
- No Cluster Management  
No infrastructure to manage, portable to Kubernetes
- Deploy source, Dockerfile or image  
Uses open source buildbacks and Cloud Build, integrates with Artifact Registry or Docker Hub

## Pricing

### Perpetual Free Tier

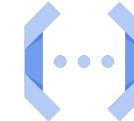
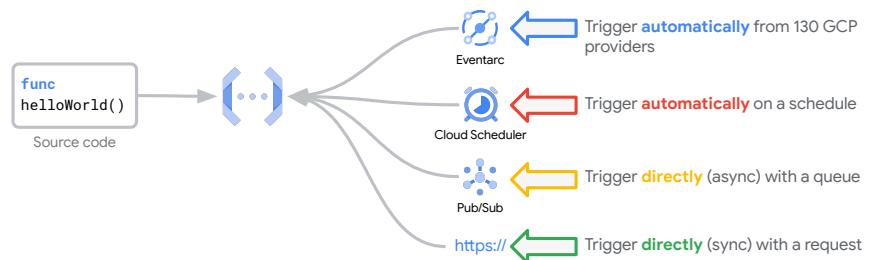
2M requests, 180K vCPU-seconds  
Per month

Pay per second of execution time  
Price varies based on memory & cpu allocation. HTTP included with request fee. Request fees waived for Always-on CPU

Google Cloud



# Cloud Functions



## Event-based applications, data transformations, ops automation

- Deploy source code  
Online editor, from bucket, using Github Actions
- Bind to trigger  
Automatically on event, or directly sync & async
- Auto-scales down to zero  
No pre-provisioning, only pay for what you use
- Supports all popular languages  
Node.js, Python, Go, Java, PHP, Ruby, .NET

## Pricing

**Perpetual Free Tier**  
2M requests, 180K vCPU-seconds  
Per month

Pay per second of execution time  
Price varies based on function size  
HTTP included with request fee

Google Cloud

# Fill in the blank



Compute  
Engine

\_\_\_\_\_ on Google's infrastructure  
Web apps, databases, machine learning, etc



Cloud  
Run

\_\_\_\_\_ containers managed by Google  
Web apps, APIs, microservices



Kubernetes  
Engine

\_\_\_\_\_ management service from Google

Web apps, databases, machine learning, etc



Cloud  
Functions

Serverless \_\_\_\_\_ managed by Google

Event-based applications, data transformations,  
ops automation

Google Cloud



# Fill in the blank



Compute  
Engine

Virtual machines on Google's infrastructure  
Web apps, databases, machine learning, etc



Cloud  
Run

Serverless containers managed by Google  
Web apps, APIs, microservices



Kubernetes  
Engine

Container management service from Google  
Web apps, databases, machine learning, etc



Cloud  
Functions

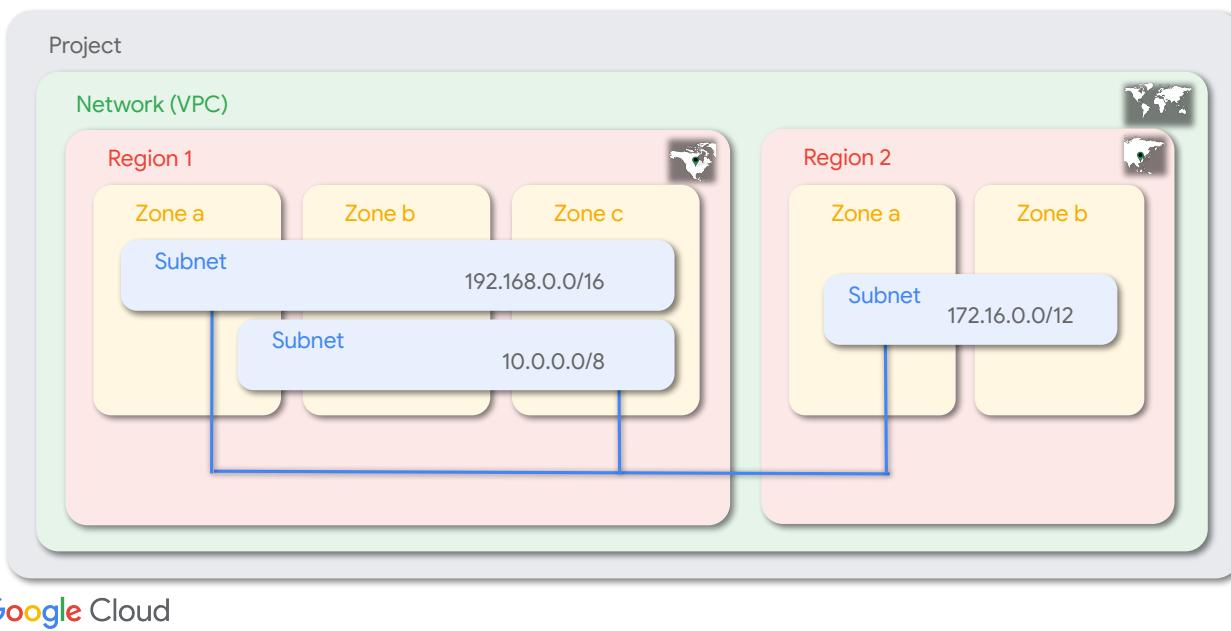
Serverless functions managed by Google  
Event-based applications, data transformations,  
ops automation

Google Cloud

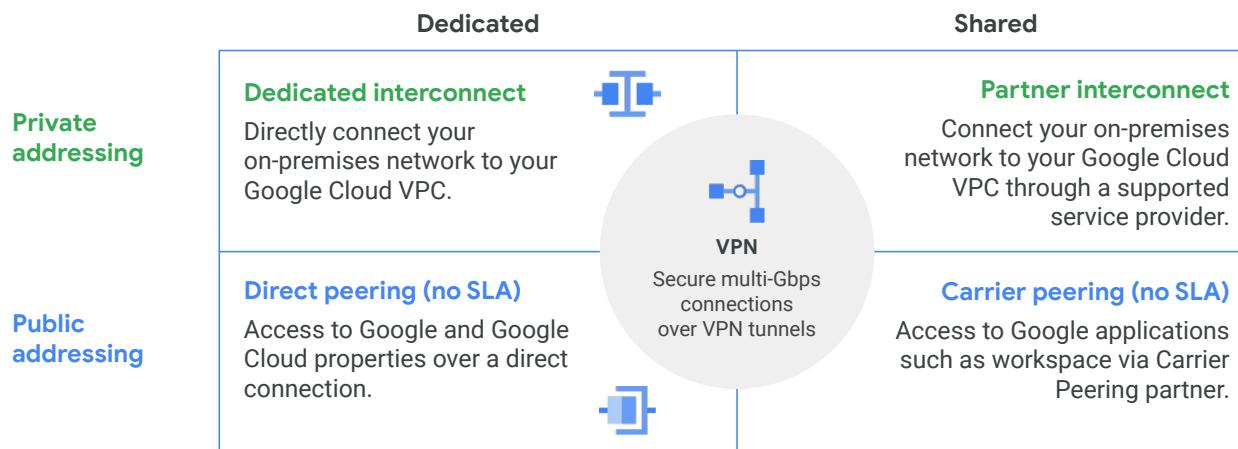


## Step 2: Connect your platform

# Google VPC



# Hybrid and multi-cloud connections



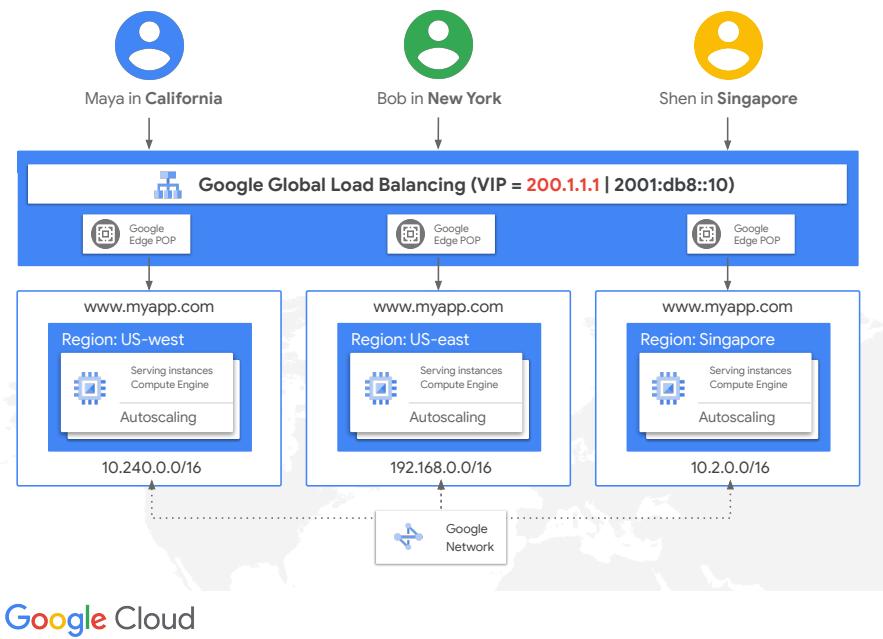
Google Cloud

## Load balancer types

	Type	Geographical scope	Network tiers	Proxy/pass-through
Internal	TCP/UDP	Regional	Premium	Pass-through
	HTTP(s)			Proxy
External	TCP/UDP	Regional	Standard / Premium	Pass-through
	HTTP(s)	Regional / Global depending on network tier	Standard / Premium	Proxy
	TCP Proxy			
	SSL Proxy			

 Google Cloud

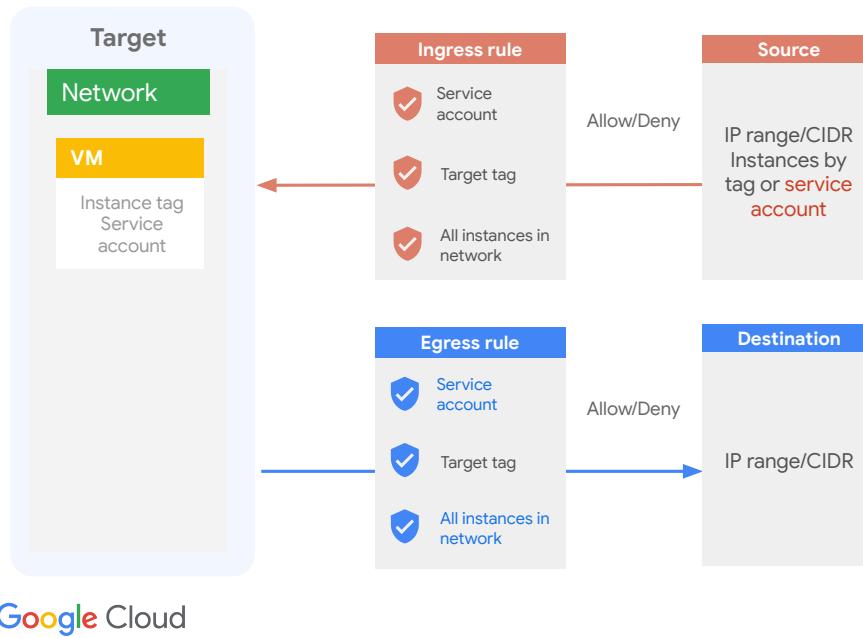
## Serving Application Traffic: External global HTTP(s) load balancing



- Global anycast IP address
- Managed service on Google Front Ends (GFE)
- Multi-regional load balancing
- Balancing mode: RATE, UTILIZATION
- Session affinity: Client IP, generated cookie, header field, & HTTP Cookie
- WebSocket support
- Terminates HTTPS traffic as close as possible to users.
- Backend traffic encryption supported.
- DDOS protection (L4)
- Integration with Cloud CDN and Cloud Armor (WAF)

## Step 3: Secure and monitor

# VPC firewall



## VPC firewall

- **Stateful** with connection tracking
- **Distributed:** enforced on underlying host

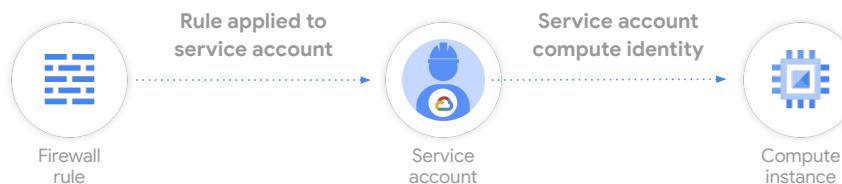
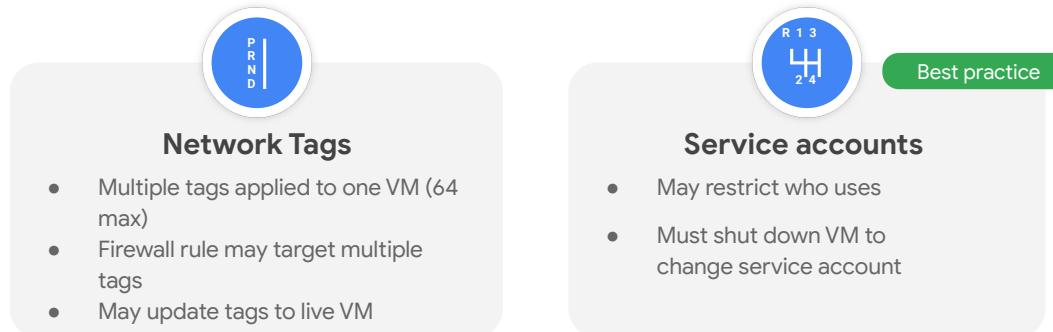
## Control paths

- VM <-> VM
- VM <-> Internet
- VM <-> On-prem

## Implied rules

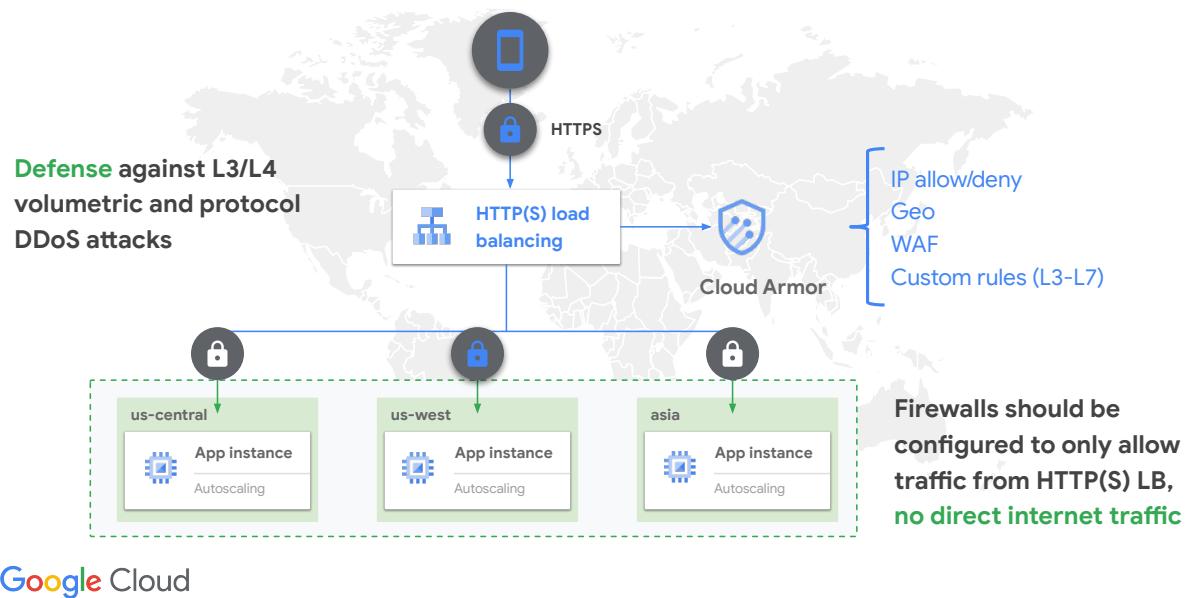
- Ingress deny
- Egress allow

## Attaching firewall rules to VMs



Google Cloud

## Cloud Armor: DDoS protection and WAF



# Cloud Operations Suite



## Logging and Error Reporting

### Cloud Logging

- Fully managed, real-time log management with storage, search, analysis and alerting at exabyte scale.

### Error reporting

- A single place to monitor error conditions from all apps and services



## Monitoring

### Cloud Monitoring

- Gain visibility into the performance, availability, and health of your applications and infrastructure.



## Application Performance Management

### Cloud Trace

- Find performance bottlenecks in production.

### Cloud Profiler

- Identify patterns of CPU, time, and memory consumption in production.

Google Cloud

# Putting it all together

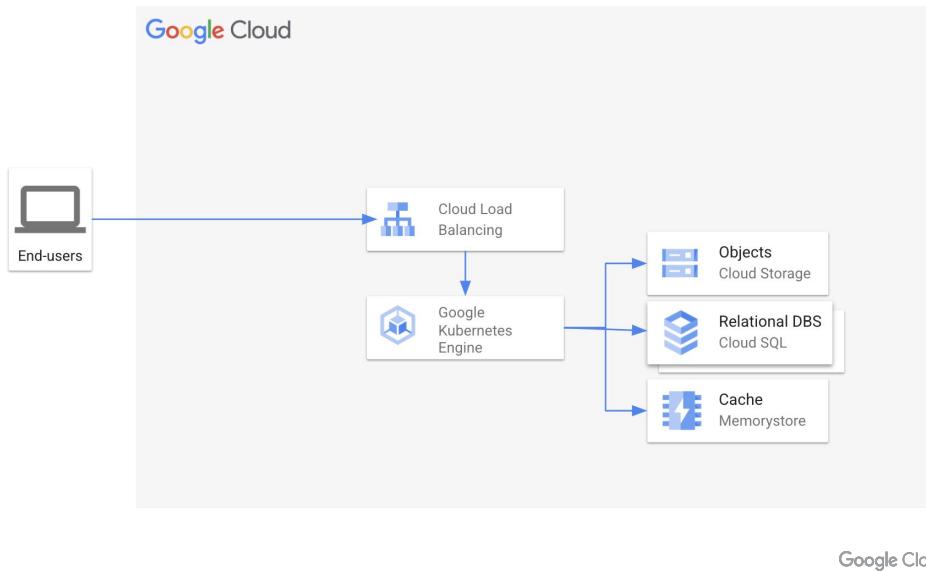
## Three-tier web application with Google Kubernetes Engine (GKE)

Example of a three-tier web application with [Google Kubernetes Engine \(GKE\)](#).

The application's services run on a [GKE](#) cluster, enabling autoscaling if necessary.

[Cloud SQL](#) is used as relational database, [Cloud Storage](#) for objects. [Memorystore](#) is used for caching to reduce access to the database on frequent queries.

Architecture: Three-tier Web Application on GKE



**Try it yourself**

**Demo :**

**[goole/loff-demo24-TTGKE](https://www.google.com/search?q=goole%2Foff+demo24-TTGKE)**

**More solutions:**

**[goole/loff-demo24-click-to-deploy](https://www.google.com/search?q=goole%2Foff+demo24-click-to-deploy)**

Google Cloud

# **What about data?**

Google Cloud

## **How to store and manage your data**

Google Cloud



**What types of data are  
you working with?**

# Managed database options

 Cloud SQL	 AlloyDB	 Spanner	 Bigtable	 Firestore
Enterprise-ready fully managed relational database service for PostgreSQL, MySQL, SQL Server	PostgreSQL-compatible database ready for enterprise level workloads	Global scale and 99.999% availability with PostgreSQL interface	Key-value database with flexible schema, single-digit millisecond low latency, and high throughput	Serverless document database with a rich development ecosystem and backend as a service
<b>Modernize with fully managed databases</b>	<b>High performance with open source compatibility and the best of Google</b>	<b>Scalability and availability for the demanding workloads</b>	<b>Fast reads and writes for high throughput workloads</b>	<b>Native integration with Google Cloud and Firebase</b>

Google Cloud



**Cymbal  
Bank**



Google Cloud

# VP of IT

Adityah



## Barriers to overcome



Operational complexity



Lack of monitoring & troubleshooting



Security & cost risks



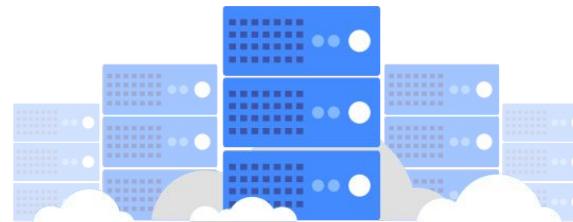
Complicated migrations

Google Cloud

# Cloud SQL

Fully managed relational database service

Support for



Google Cloud



1  
Fully managed  
Automation

2  
Reduced risk  
Embedded Tooling

3  
Better observability  
Query Insights

4  
Easy migrations  
Database Migration Service

5  
Higher performance  
PostgreSQL Workloads

6  
HTAP workloads  
AI & ML Innovation

Google Cloud   AlloyDB Golden Demo Test   Search Products, resources, docs (/)

SQL Instances + CREATE INSTANCE MIGRATE DATA HELP ASSISTANT

## Automated and fully managed

Provisioning



### Cloud SQL Instances

Cloud SQL instances are fully managed, relational MySQL, PostgreSQL, and SQL Server databases. Google handles replication, patch management, and database management to ensure availability and performance. [Learn more](#)

To get started with Cloud SQL, you can create a new instance or use Database Migration Service to migrate your SQL database to Google Cloud.

[CREATE INSTANCE](#) [MIGRATE DATA](#)



1  
Fully managed  
Automation

2  
Reduced risk  
Embedded Tooling

3  
Better observability  
Query Insights

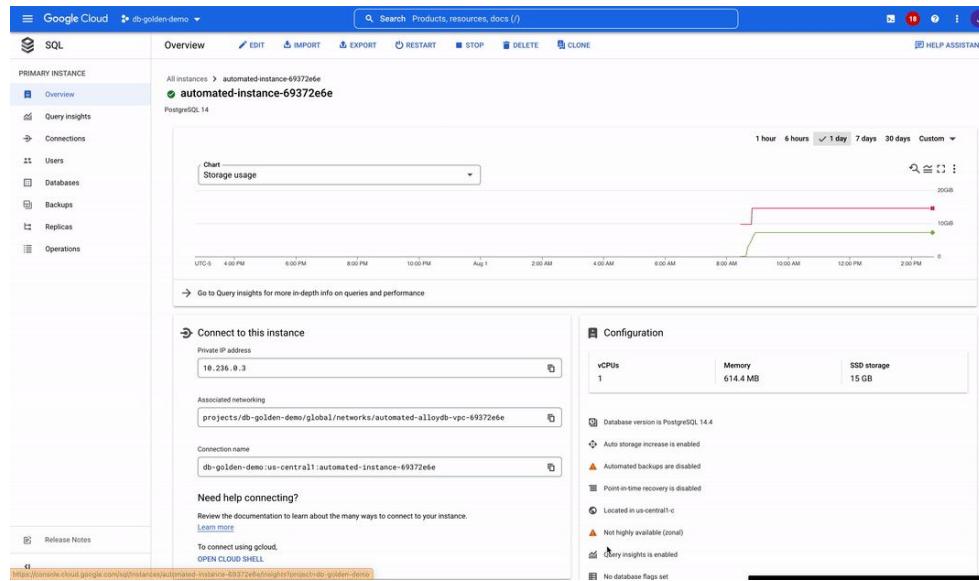
4  
Easy migrations  
Database Migration Service

5  
Higher performance  
PostgreSQL Workloads

6  
HTAP workloads  
AI & ML Innovation

## Automated and fully managed

Auto storage increase



Google Cloud



1  
Fully managed  
Automation

2  
Reduced risk  
Embedded Tooling

3  
Better observability  
Query Insights

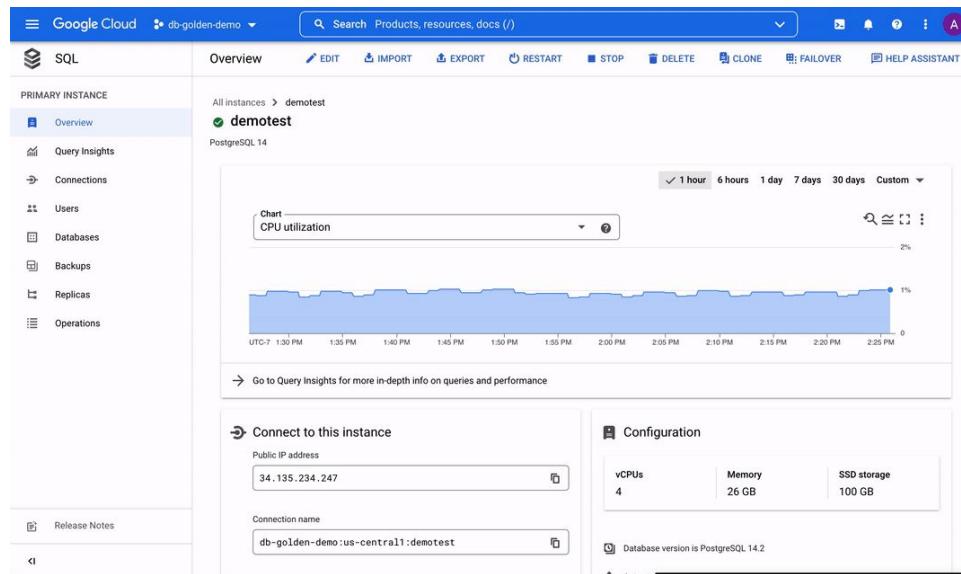
4  
Easy migrations  
Database Migration Service

5  
Higher performance  
PostgreSQL Workloads

6  
HTAP workloads  
AI & ML Innovation

## Automated and fully managed

Connectivity



Google Cloud



1  
Fully managed  
Automation

2  
Reduced risk  
Embedded Tooling

3  
Better observability  
Query Insights

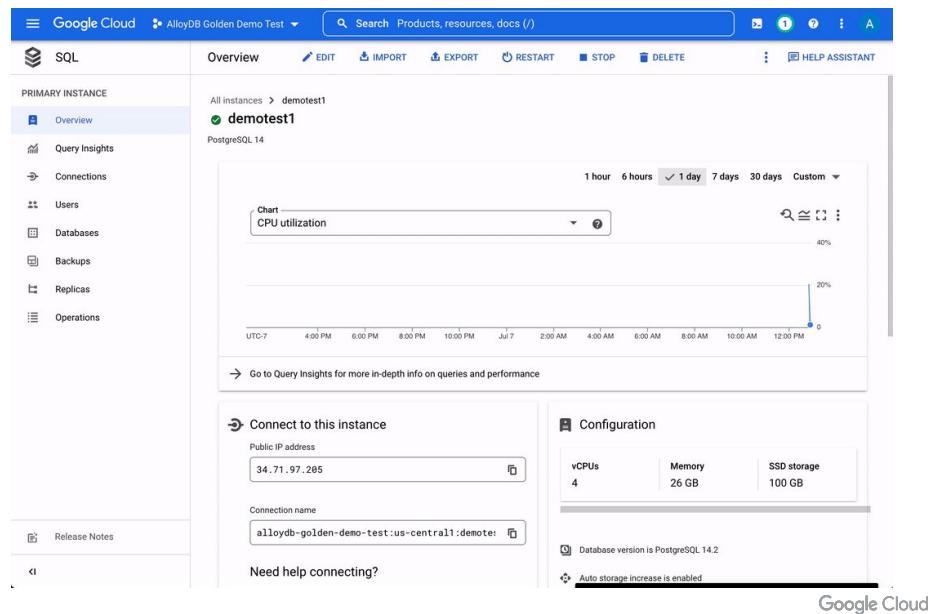
4  
Easy migrations  
Database Migration Service

5  
Higher performance  
PostgreSQL Workloads

6  
HTAP workloads  
AI & ML Innovation

## Automated and fully managed

Backup and recovery





1  
Fully managed  
Automation

2  
Reduced risk  
Embedded Tooling

3  
Better observability  
Query Insights

4  
Easy migrations  
Database Migration Service

5  
Higher performance  
PostgreSQL Workloads

6  
HTAP workloads  
AI & ML Innovation

## Automated and fully managed

Maintenance

The screenshot shows the Google Cloud SQL interface for managing a PostgreSQL database instance named 'db-golden-demo'. The left sidebar lists primary instance management options like Overview, Query Insights, Connections, Users, Databases, Backups, Replicas, and Operations. The main content area displays the 'Edit demotest' configuration page under the 'SPECIFY ZONES' tab. It includes sections for 'Customize your instance' (Machine type, Storage, Connections, Backups, Maintenance), 'Flags' (No flags set), and 'Insights' (Insights enabled). On the right, a 'Summary' table provides detailed technical specifications:

Region	us-central1 (Iowa)
DB Version	PostgreSQL 14.2
vCPU	4 vCPU
Memory	26 GB
Storage	100 GB
Network throughput (MB/s)	1,000 of 2,000
Disk throughput (MB/s)	Read: 48.0 of 240.0 Write: 48.0 of 240.0
IOPS	Read: 3,000 of 15,000 Write: 3,000 of 15,000
Connections	Public IP
Backup	Automated
Availability	Multiple zones (Highly available)
Point-in-time recovery	Enabled

Google Cloud



① Fully managed Automation

② Reduced risk Embedded Tooling

③ Better observability Query Insights

④ Easy migrations Database Migration Service

⑤ Higher performance PostgreSQL Workloads

⑥ HTAP workloads AI & ML Innovation

## Risk Reduction

### Security

The screenshot shows the Google Cloud IAM & Admin interface. The left sidebar includes options like IAM, Identity & Organization, Policy Troubleshooter, Policy Analyzer, Organization Policies, Service Accounts, Workload Identity Federation, Labels, Tags, Settings, Privacy & Security, Identity-Aware Proxy, Roles, Audit Logs, Asset Inventory, Manage Resources, and Release Notes. The main area is titled "Permissions for project 'db-golden-demo'" and lists principals and their roles. A table provides details for each principal:

Type	Principal	Name	Role	Security insights	Inheritance
User	195449165820@cloudfs.gserviceaccount.com	Google APIs Service Agent	Editor	0/100 excess permissions	✓
User	admin@1987984870407.aitostrat.com	Tastir Imam	Owner	5647/5674 excess permissions	✓
Group	css_database_demo_group@1987984870407.aitostrat.com	BigQuery Admin Cloud AlloyDB Admin Cloud SQL Admin Create Service Accounts Editor Organization Administrator Vertex AI Administrator	103/104 excess permissions 28/29 excess permissions 70/71 excess permissions 3/5 excess permissions 5211/5228 excess permissions 221/222 excess permissions	0/100 excess permissions 0/100 excess permissions 0/100 excess permissions 0/100 excess permissions 0/100 excess permissions 0/100 excess permissions	✓

Google Cloud



1  
Fully managed Automation

2  
Reduced risk  
Embedded Tooling

3  
Better observability  
Query Insights

4  
Easy migrations  
Database Migration Service

5  
Higher performance  
PostgreSQL Workloads

6  
HTAP workloads  
AI & ML Innovation

## Risk Reduction

### Cost Recommender

The screenshot shows the Google Cloud Cost Recommender interface under the 'RECOMMENDATIONS' tab. It displays six recommendations:

- Reduce Cloud SQL instance cost**: Stop idle and resize underutilized Cloud SQL resources to reduce costs. Cost savings: \$1,046.12/month estimate. Action: Shut down instance to save \$195.47/month + 2 more. View all.
- Prevent downtime for Cloud SQL instances**: Prevent Cloud SQL instances from running out of disk space and going offline by automatically creating storage. Performance: 1 instances at or nearing capacity. Action: Enable automatic storage increase for aar.ukruk. View all.
- Unused Compute Engine resources**: Back up and delete unused resources to reduce costs. Cost savings: \$202.13/month estimate. Action: Delete disk to save \$0.68/month + 35 more. View all.
- Limit cross-project impersonations**: Increase the security of your cloud by limiting access from outside project impersonation capabilities. Security: 17,462 excess permissions estimate. Action: Remove editor role recommendation + 5 more. View all.
- Change project-level IAM role grants**: Increase the security of your cloud by removing or replacing overly permissive roles granted to principals at the project level. Security: 942,032 excess permissions estimate. Action: Remove owner role recommendation + 456 more. View all.
- Secure Cloud Run services**: Take recommended action to increase the security of your Cloud Run service. Security: Increased security. Action: Create new service account + 2 more. View all.

Google Cloud



① Fully managed Automation

② Reduced risk Embedded Tooling

③ Better observability  
Query Insights

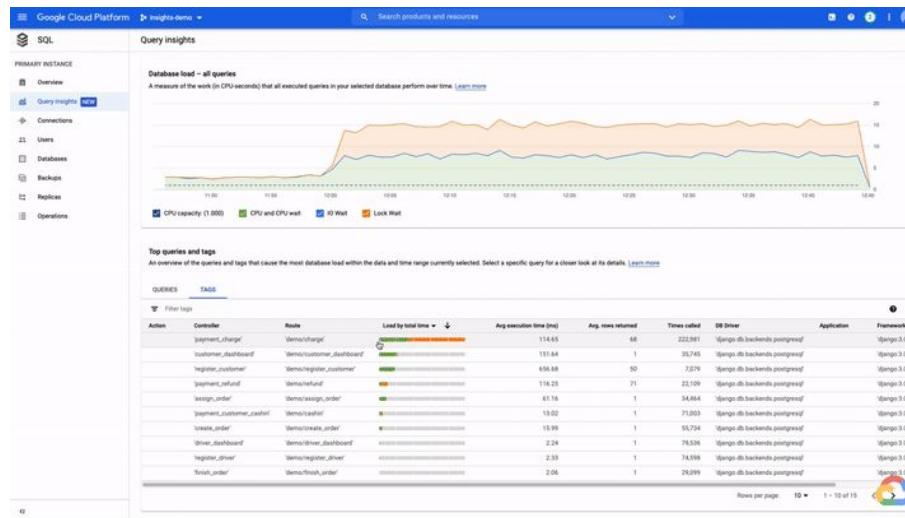
④ Easy migrations  
Database Migration Service

⑤ Higher performance  
PostgreSQL Workloads

⑥ HTAP workloads  
AI & ML Innovation

## Productivity and observability

Cloud SQL Query Insights tool provides monitoring and diagnostics that helps to detect and fix query performance problems.



Google Cloud



1 Fully managed Automation

2 Reduced risk Embedded Tooling

3 Better observability Query Insights

4 Easy migrations

Easy migrations

5 Higher performance Database Migration Service PostgreSQL Workloads

6 HTAP workloads AI & ML Innovation

## Easy Migration

Migration to Cloud SQL using Database Migration Service is serverless, and easy to use.

The screenshot shows the 'Create a migration job' wizard in Google Cloud. The current step is 'Describe your migration job'. The left sidebar lists six steps: 1. Get started (Not configured), 2. Define a source (Not configured), 3. Create a destination (Not configured), 4. Define connectivity method (Not configured), 5. Test and create migration job (Not tested). Step 4 is highlighted in green. The main form fields include:

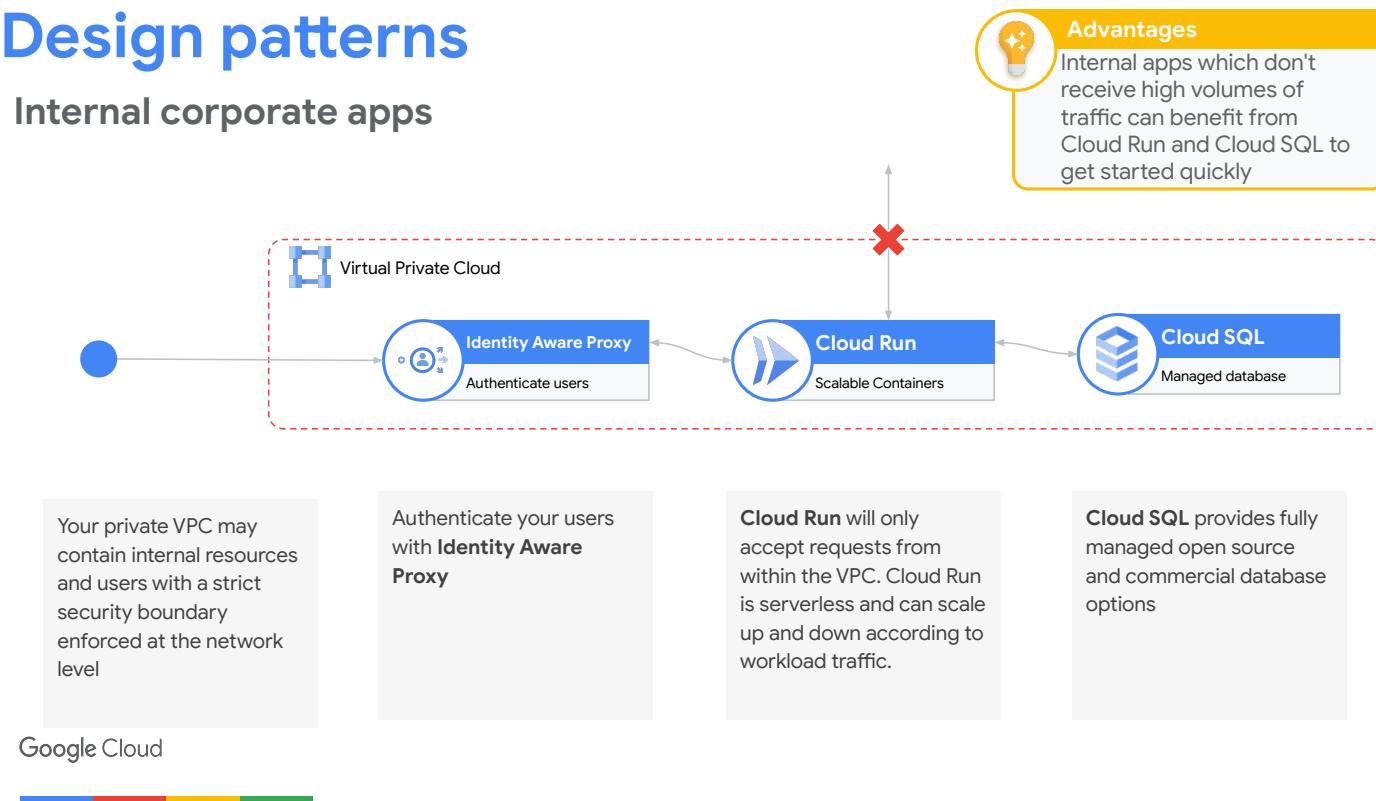
- Migration job name:** postgresql-on-premises-cloudsql-migration
- Migration job ID:** postgresql-on-premises-cloudsql-migration
- Source database engine:** (dropdown)
- Destination region:** us-central1 (Iowa)
- Migration job type:** (dropdown)

Below the form, a note says: "Before you continue, review the prerequisites".

Google Cloud

# Design patterns

## Internal corporate apps



# VP of IT

Adityah



## Cloud SQL



Fully managed operations



Risk reduction and governance



Better observability



Easy migration

Google Cloud

## Cymbal Bank's new initiative

Offering customers a unified experience for both banking and investment account management



Google Cloud

# Head of Digital Banking

Elena



Increasing trade volumes and volatility



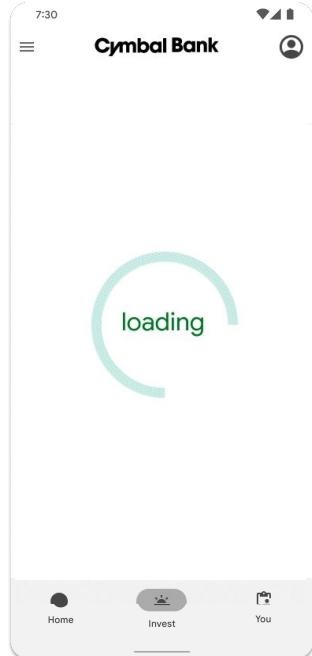
Need real-time insights for my customers

Google Cloud



## Real-time Operational Insights

AlloyDB vs. Conventional PostgreSQL:  
A side-by-side comparison



Status quo



w/ AlloyDB

Google Cloud

# AlloyDB



**2x faster**

than Amazon's comparable  
PostgreSQL-compatible service  
for transactional workloads

**100x faster**

for analytical queries than standard  
PostgreSQL

Google Cloud



① Fully managed Automation

② Reduced risk Embedded Tooling

③ Better observability Query Insights

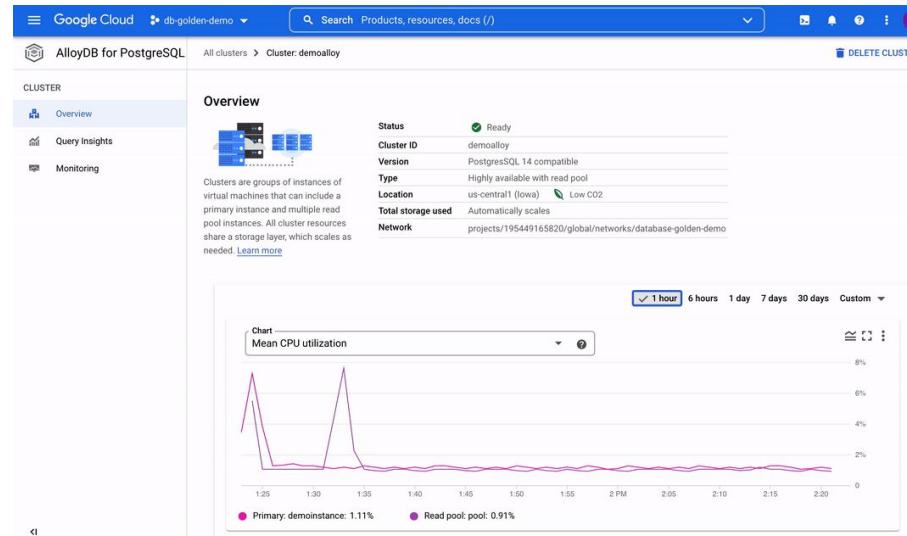
④ Easy migrations Database Migration Service

⑤ Higher performance PostgreSQL Workloads

⑥ HTAP workloads AI & ML Innovation

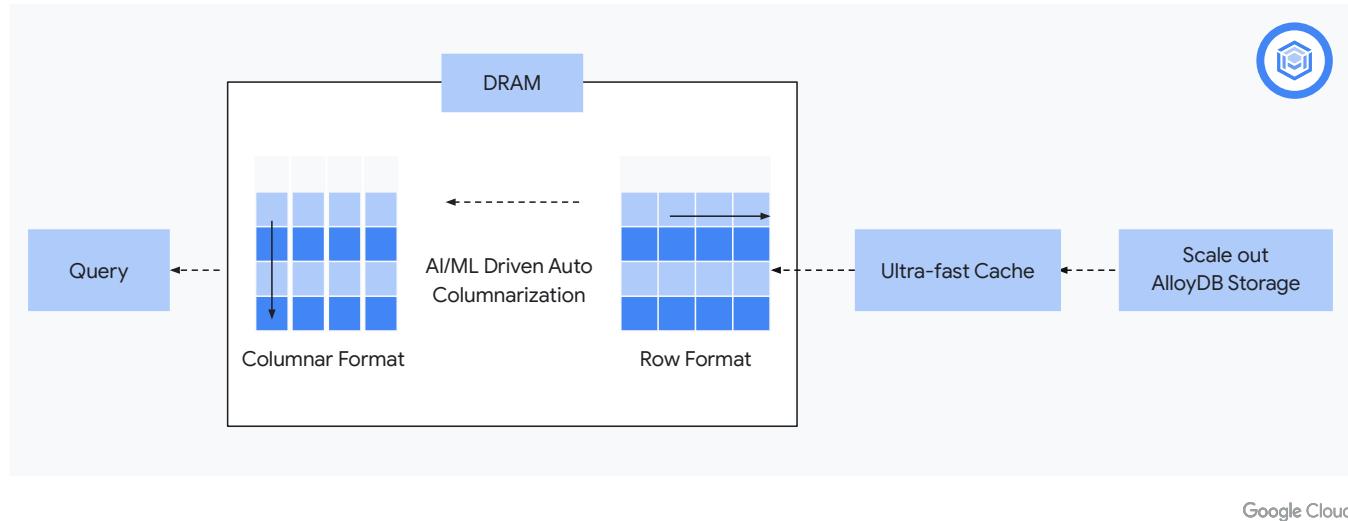
## Real Time Operational Analytics

Hybrid transactional and analytical processing (HTAP)



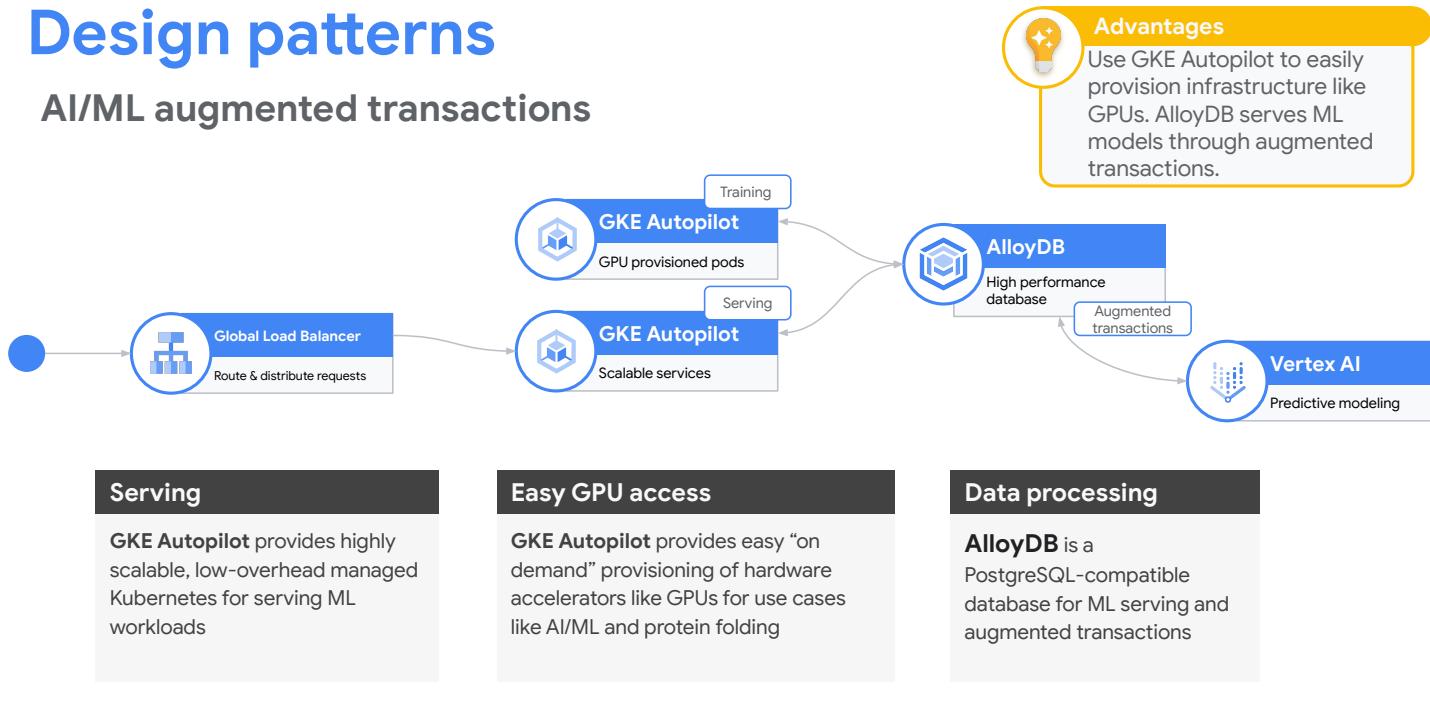
Google Cloud

The **columnar cache** that's powering the dashboard is itself powered by ML based on the unique characteristics of the workload



# Design patterns

## AI/ML augmented transactions



Google Cloud

# Head of Digital Banking

Elena



Trade volume scaling without issues

New data value created for  
Cymbal Bank and their customers

Predictable and transparent pricing

Google Cloud

# Head of Digital Banking: Initiative #2

Elena

## A Specific Technical Challenge:

-  Digital trading app uses HBase for delivering real time recommendations to the user and for detecting fraud
-  HBase is not scaling to the needs of the app - outages occur during heavy trading periods
-  HBase is costly and difficult to maintain



Google Cloud



## Why modernize with Bigtable?

- Bigtable is a managed NoSQL service that supports the open HBase API - [no code lock in](#)
- We provide the tools to help automate assessment and subsequent migration [easier and faster](#)
- We have migrated [hundreds of customers](#) across every industry
- [Google uses Bigtable](#) extensively in its most popular apps



## Bigtable Common Use Cases



Personalization  
and  
Recommendations  
with real-time  
clicks and  
transactions



Fraud Detection  
with real-time  
events



Predictive  
Maintenance with  
IoT / Machine Data

Metadata  
Management for  
customer and  
product metadata

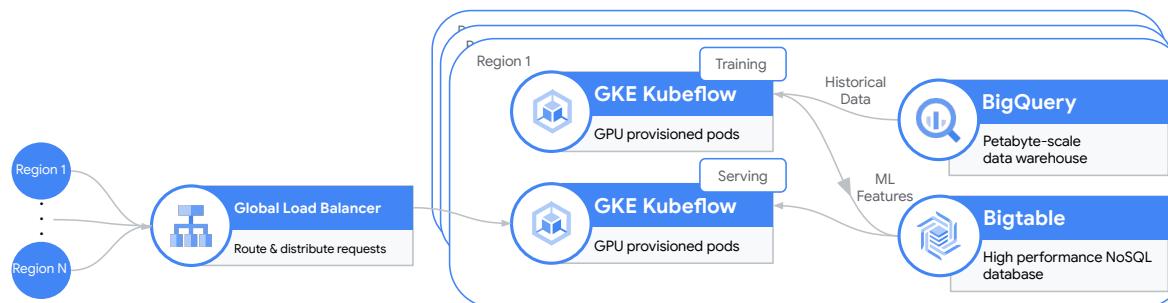


Google Cloud

# Design patterns

## AI/ML training and low-latency serving

**Advantages**  
Use GKE and Bigtable for workloads that require low-latency and high throughput



ML workflow	Historical analysis	Online inference & model serving	Global deployments
Kubeflow on GKE simplifies the build-train-deploy lifecycle.	BigQuery provides the historical data for analysis, feature engineering, offline model training and evaluation.	Bigtable provides high read & write throughput with low latency. It can adapt to fluctuating demand with built-in autoscaling. Couple with MLRun, Tensorflow Serving or KServe on GKE, to serve real-time predictions at scale for ML framework of your choice.	Bigtable supports read/write clusters in <b>up to 8 regions of your choice</b> for fast regional access, higher resiliency and data residency compliance.

Google Cloud

# Head of Digital Banking

Elena



On premise HBase replaced with fully managed Bigtable



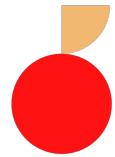
Automatic scaling up and down to meet the needs of the app - no more capacity related outages



Bigtable

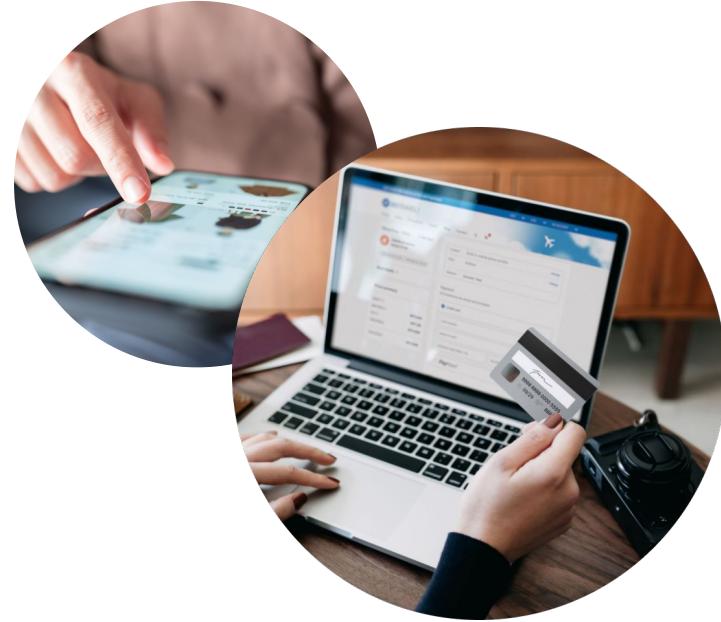


Google Cloud



# Cymbol Superstore

Retailer with a large,  
rapidly growing digital business



Google Cloud

# Director of Ecommerce

Elma



Proprietary + Confidential



Online store application has outages during heavy use periods - no managed database product scales well enough



Application must run in multiple regions, but it must look like it is one database - we must have a consistent view of all data



Even an outage for a few minutes (planned or unplanned) is too expensive  
- want 99.999% SLA

Google Cloud



# Cloud Spanner

Not just another \_\_\_\_ database.



Performance

and



Scale

and



SQL

and



Fully Managed Service

Google Cloud



## What is Unique About Spanner?

Battle tested by Google



### Relational semantics

Schemas, ACID  
transactions, SQL



### Horizontal scale

99.999% SLA, fully  
managed, and scalable

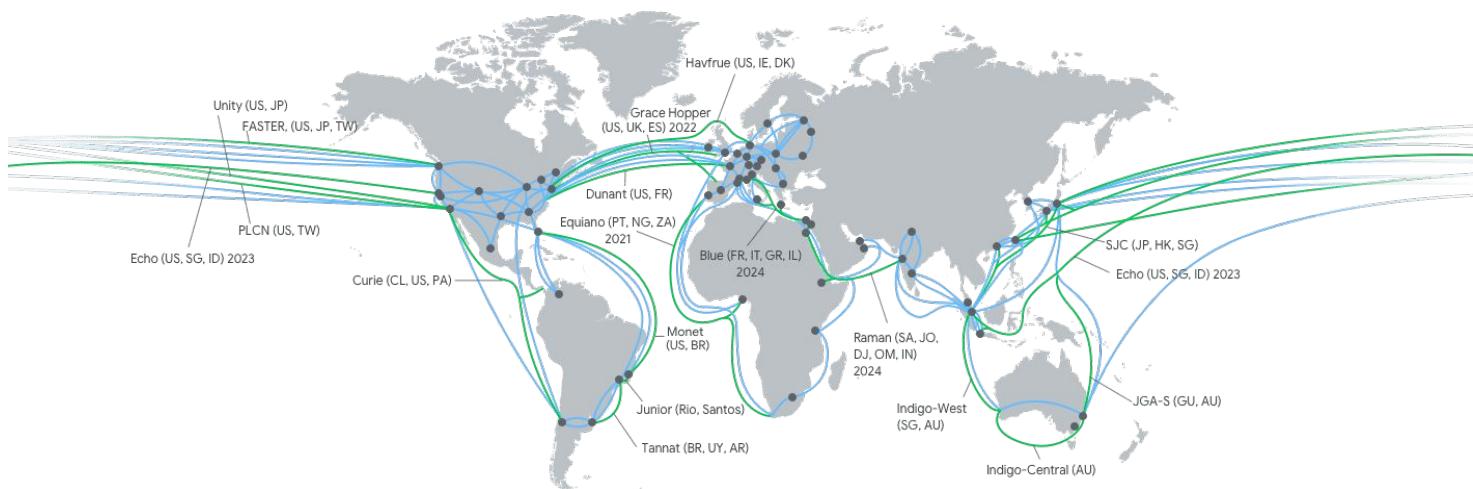
*Spanner is a fully managed database service that processes over 2 billion requests per second at peak*

Google Cloud



# Google's Unique Worldwide Network

Powered by Google

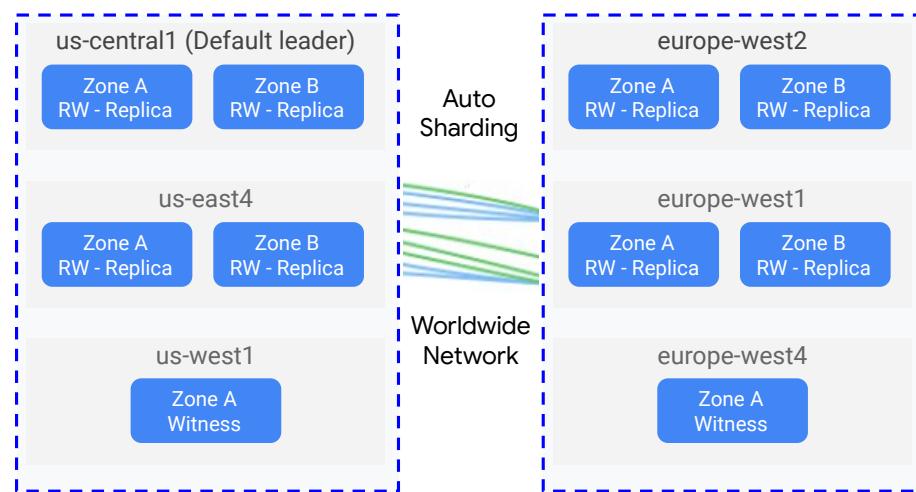


Google Cloud



# Availability at any scale

- Industry leading 99.999% SLA  
RPO = 0, RTO = 0
  - Multiple read-write regions  
with a consistent, single  
view of all data
  - No maintenance window
  - Online schema changes



## Multi region example





Cymbal  
Superstore

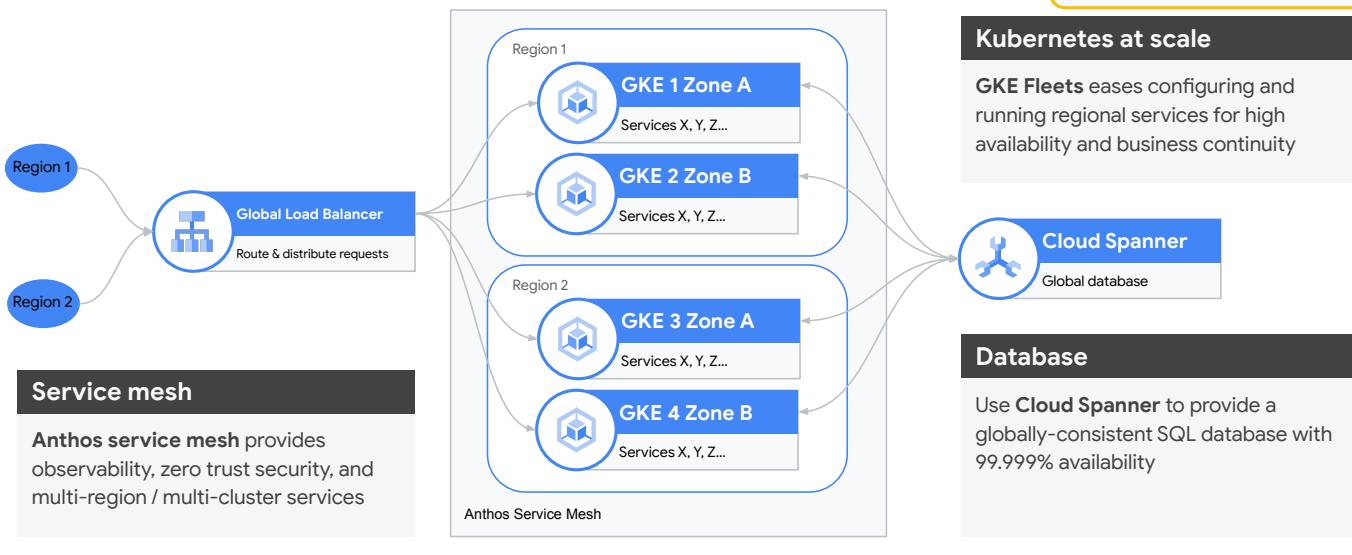
## Spanner Continued Innovation

Democratized Access	Enterprise Readiness	Developer Experience	Performance and Flexibility
<p>Reducing the barriers for teams to access Spanner's unique benefits:</p> <ul style="list-style-type: none"><li>• Lower entry price: Granular instances</li><li>• Committed use discounts: Up to 40%</li><li>• Familiarity, portability: PostgreSQL interface</li><li>• Free Trial Instance</li></ul> 	<p>Powering mission-critical applications in regulated industries:</p> <ul style="list-style-type: none"><li>• Point-in-time recovery</li><li>• CMEK organization policy</li><li>• Fine-grained access control</li></ul> <small>PREVIEW</small> 	<p>Simplifying integration into applications and developer workflows:</p> <ul style="list-style-type: none"><li>• Query Insights</li><li>• BigQuery federation</li><li>• Cost-based query optimizer</li><li>• Self-service optimizer stats updates (ANALYZE)</li><li>• Open Telemetry, Grafana support</li></ul> 	<p>Realizing the benefits of the cloud for operational workloads:</p> <ul style="list-style-type: none"><li>• Doubled storage: 4TB/node</li><li>• Mutation limit: 40k</li><li>• Leader placement</li><li>• New multi-regions</li></ul> 

Google Cloud

# Design patterns

## Business continuity for enterprise apps



Google Cloud

## Firestore's Mission

Unlock application innovation with **simplicity**, **speed**, and **confidence**.

### What makes Firestore special?

#### Serverless document database

Built for developers: scalable, schemaless (e.g. JSON), ACID transactions, pay only for what you use

#### Backend as a service

Connect directly, no middle tier needed, development to production in days

#### Built-in real-time sync & offline mode

Easily build rich apps with live sync, that also fallback to offline caching when a client device loses network connectivity

Firestore has more than **250,000** monthly active developers

Google Cloud

## Differentiated infrastructure

**01 Serverless:** Google battle-tested, automatic, elastic scaling, that helps minimize capacity planning. You pay for what you use.

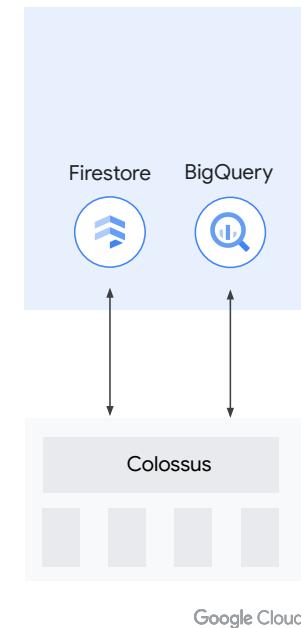
**02 Serializable transactions:** serializable ACID transactions with up to 500 documents.

**03 High availability:** up to 99.999% SLA availability (multi-region), 0 RPO/RTO for zonal and regional failures, and no maintenance window downtimes.

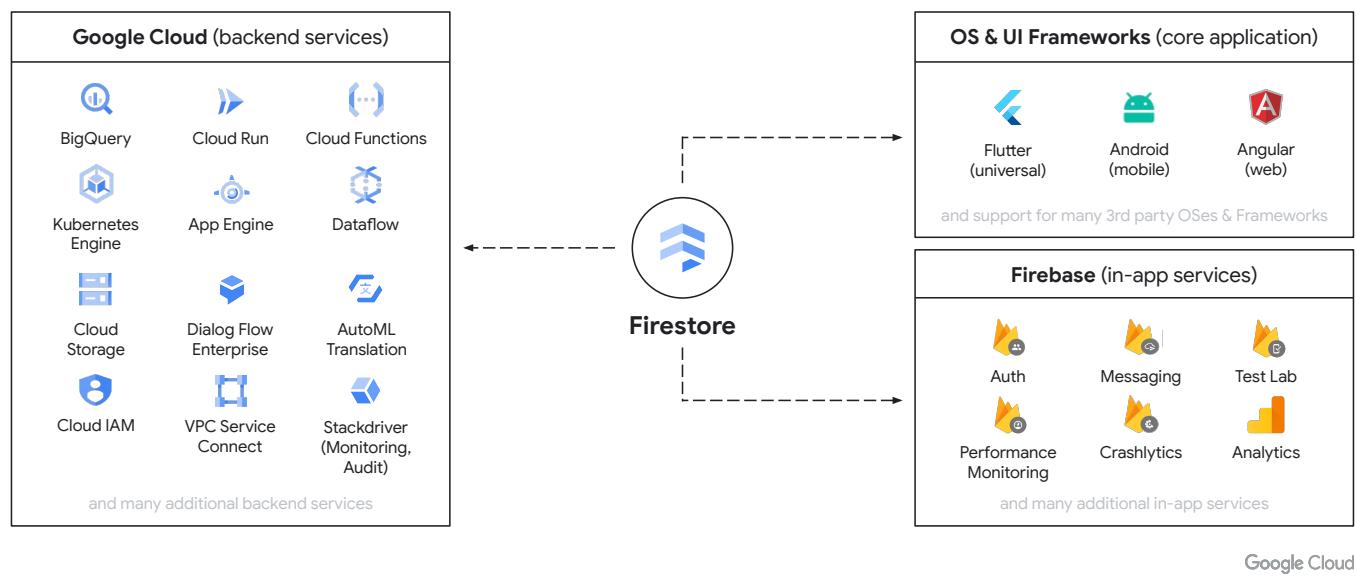
**04 External consistency:** external consistency based on TrueTime.

**05 Multi-version concurrency control:** allows customers to perform queries both on most recent and past data.

**06 Shared Colossus storage:** performant data access for load, storage and retrieval for a variety of use cases - OLTP, analytics, AI/ML



## Integrated full-stack application development ecosystem from Google



# Director of Ecommerce

Elma



Proprietary + Confidential



Massively scaling application even during peak times of heavy load. No outages due to capacity



Data resides in multiple regions, but consistent queries and reports look like it is one database



No data loss - no downtime for maintenance - instant recovery from failures

Google Cloud

# Managed database options

				
Enterprise-ready fully managed relational database service for PostgreSQL, MySQL, SQL Server	PostgreSQL-compatible database ready for enterprise level workloads	Global scale and 99.999% availability with PostgreSQL interface	Key-value database with flexible schema, single-digit millisecond low latency, and high throughput	Serverless document database with a rich development ecosystem and backend as a service
<b>Modernize with fully managed databases</b>	<b>High performance with open source compatibility and the best of Google</b>	<b>Scalability and availability for the demanding workloads</b>	<b>Fast reads and writes for high throughput workloads</b>	<b>Native integration with Google Cloud and Firebase</b>

Google Cloud

# Managed database options

 Cloud SQL	 AlloyDB	 Spanner	 Bigtable	 Firestore
Enterprise-ready fully managed relational database service for PostgreSQL, MySQL, SQL Server	PostgreSQL-compatible database ready for enterprise level workloads	Global scale and 99.999% availability with PostgreSQL interface	Key-value database with flexible schema, single-digit millisecond low latency, and high throughput	Serverless document database with a rich development ecosystem and backend as a service
<b>Modernize with fully managed databases</b>	<b>High performance with open source compatibility and the best of Google</b>	<b>Scalability and availability for the demanding workloads</b>	<b>Fast reads and writes for high throughput workloads</b>	<b>Native integration with Google Cloud and Firebase</b>

Google Cloud

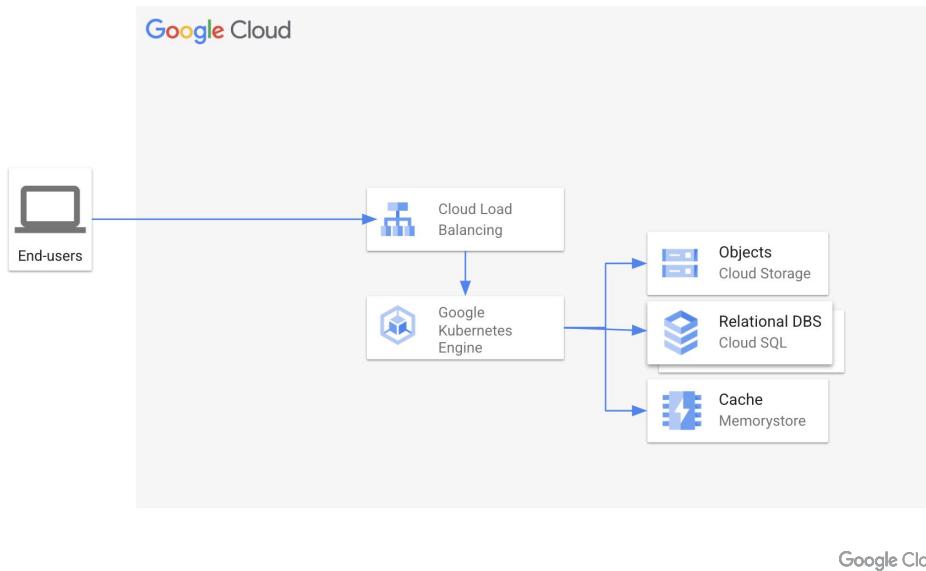
## Three-tier web application with Google Kubernetes Engine (GKE)

Example of a three-tier web application with [Google Kubernetes Engine \(GKE\)](#).

The application's services run on a [GKE](#) cluster, enabling autoscaling if necessary.

[Cloud SQL](#) is used as relational database, [Cloud Storage](#) for objects. [Memorystore](#) is used for caching to reduce access to the database on frequent queries.

Architecture: Three-tier Web Application on GKE



**Try it yourself**

**Demo :**

**[goole/loff-demo24-TTGKE](https://www.google.com/search?q=goole%2Foff+demo24-TTGKE)**

**More solutions:**

**[goole/loff-demo24-click-to-deploy](https://www.google.com/search?q=goole%2Foff+demo24-click-to-deploy)**

Google Cloud

# How to analyze your data

Google Cloud



# **What is a data warehouse?**

Google Cloud

## What is a data warehouse?

A data warehouse is an enterprise system used for the **analysis and reporting of structured and semi-structured data from multiple sources**, such as point-of-sale transactions, marketing automation, customer relationship management, and more.

Google Cloud

# Google BigQuery

Google Cloud Platform's  
**enterprise data warehouse**  
for analytics

Gigabyte- to **petabyte-scale**  
storage and SQL queries

**Encrypted**, durable,  
And highly available



Fully managed and **serverless**  
for maximum agility and scale

Unique

**Real-time** insights from streaming data

Unique

Built-in **ML** for out-of-the-box  
predictive insights

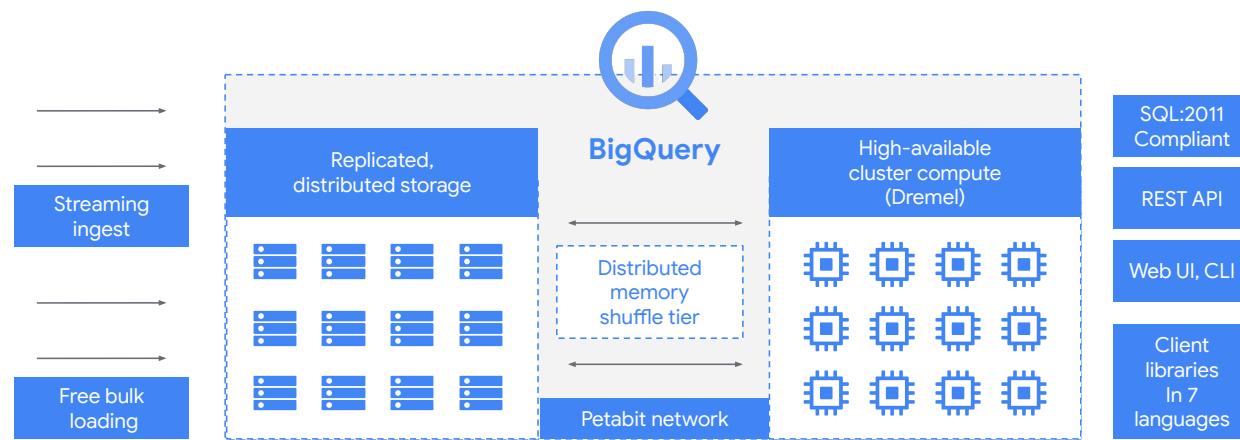
Unique

High-speed, in-memory **BI Engine**  
for faster reporting and analysis

Unique

Google Cloud

# BigQuery Architecture



Google Cloud

# Completely serverless

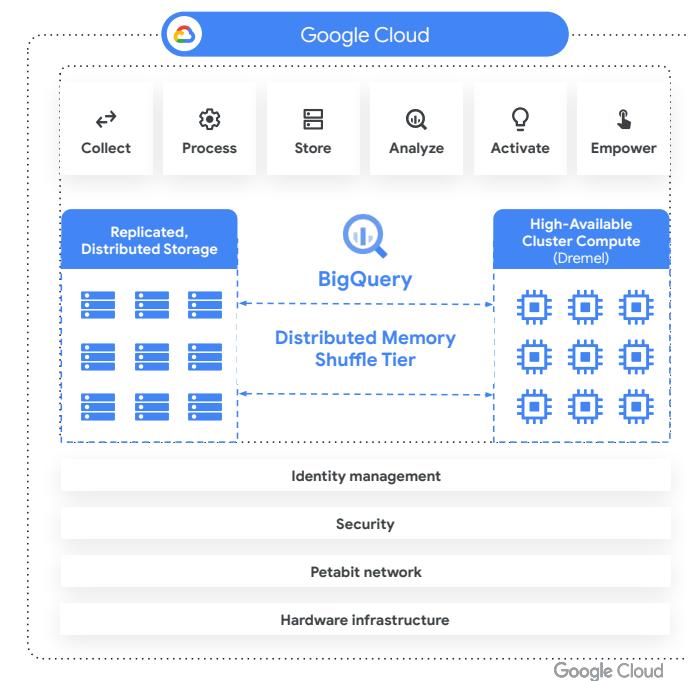
## Why BigQuery?

- Simplifies capacity management
  - Dynamically adjusts to demand
- VS
- Plan, manage, pay VMs
  - Limit use due to capacity restrictions

## Elastic Computing

Distributed storage and compute with high bandwidth including distributed petabyte scale in-memory storage for temp data and state:

- Auto-start and auto-pause
- 0-Second warm up to get maximum performance
- Accelerate queries in flight



# Stream Analytics

## Why BigQuery?

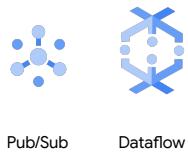
- Leverage event-driven analysis with built-in streaming capabilities

VS

- Leverage historical analysis and batch processing

## Analyze business events in real time

Move your business to event-driven action for logstream, clickstream, and sensor data to enable use cases like anomaly detection and continuous intelligence.



Pub/Sub      Dataflow

Streaming API



Stream API for immediate ingestion of events.

No additional buffer and file coalescence to manage.

BI Engine



Looker      Partner ecosystem

Google Cloud

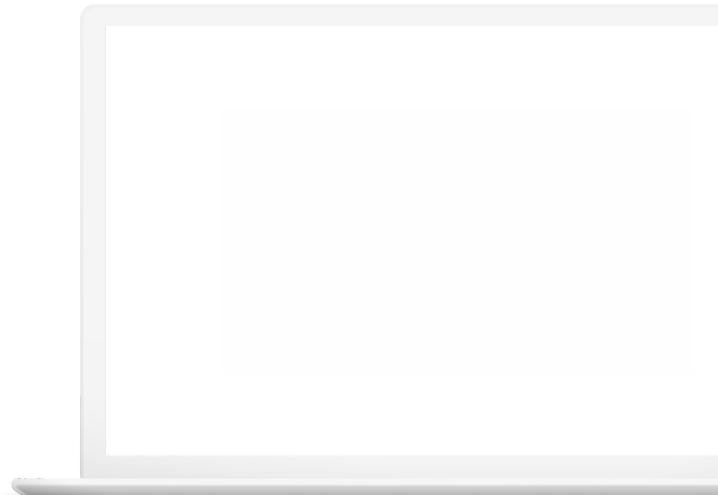
# Built-in AI/ML | BQML

## Why BigQuery?

- Provide ML access to more users through a simple SQL interface
- VS
- Require every ML use case to go through more specialized systems that require advanced skill sets

## Machine Learning for all Built-in ML with SQL

- Execute, iterate, and automate ML initiatives all within BigQuery using predefined models
- Leverage external models developed in Tensorflow directly from SQL
- Export developed models for use in Vertex AI



Google Cloud

# Built-in AI/ML | Vertex AI

## Why BigQuery?

- Simplify operationalizing models
- Manage MLOps in an environment designed to work with your Data Warehouse

VS

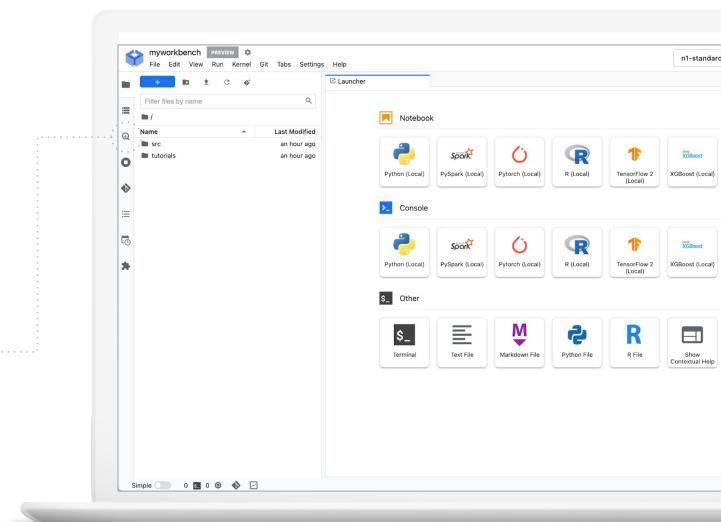
- Only use models for experiments and not operations
- Manage separate tooling for ML workloads and data integrations

## Direct integration with Vertex AI Workbench

Data scientists can use Spark, Pytorch, and Tensorflow in connection with data from BigQuery with a simple integration

## Complete MLOps environment

Vertex AI provides an environment for continuous training and deployment workflows with ML pipelines



Google Cloud

# BI Engine for BigQuery

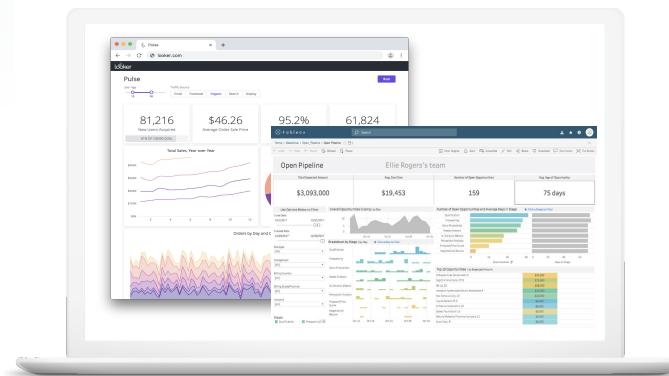
Powering real-time analytics

## Why BigQuery?

- Accelerate SQL queries from any tool with automatic caching
- Hand optimize queries or synchronize with offline BI tool caches

VS

- Adaptive cache automatically keeps hottest data in memory
- Vectorized runtime means BigQuery queries execute in sub-second
- Use BigQuery SQL directly from Tableau, Looker, PowerBI, BigQuery UI, etc.



# Data governance & security out of the box

Protecting against data exfiltration

## Why BigQuery?

- A security-first approach that brings Google's security thought leadership directly to your data
- Security as a secondary featureset

VS



Fine-grained access controls through BigQuery column-level security & integration with [Cloud IAM](#)



Metadata management and data discovery in BigQuery through [Dataplex](#)



Default encrypted with data encryption in transit and at rest



Discover, classify and redact sensitive data in GCP with [Cloud DLP](#)



Shared datasets helps securely share read-only data both internally and externally at scale

Google Cloud

**Find in the documentation:**

**Q: What is the maximum  
number of tables you can  
create in BigQuery?**

**Find in the documentation:**

**Q: What is the maximum number of tables you can create in BigQuery?**

**A: Unlimited**

<https://cloud.google.com/bigquery/quotas>

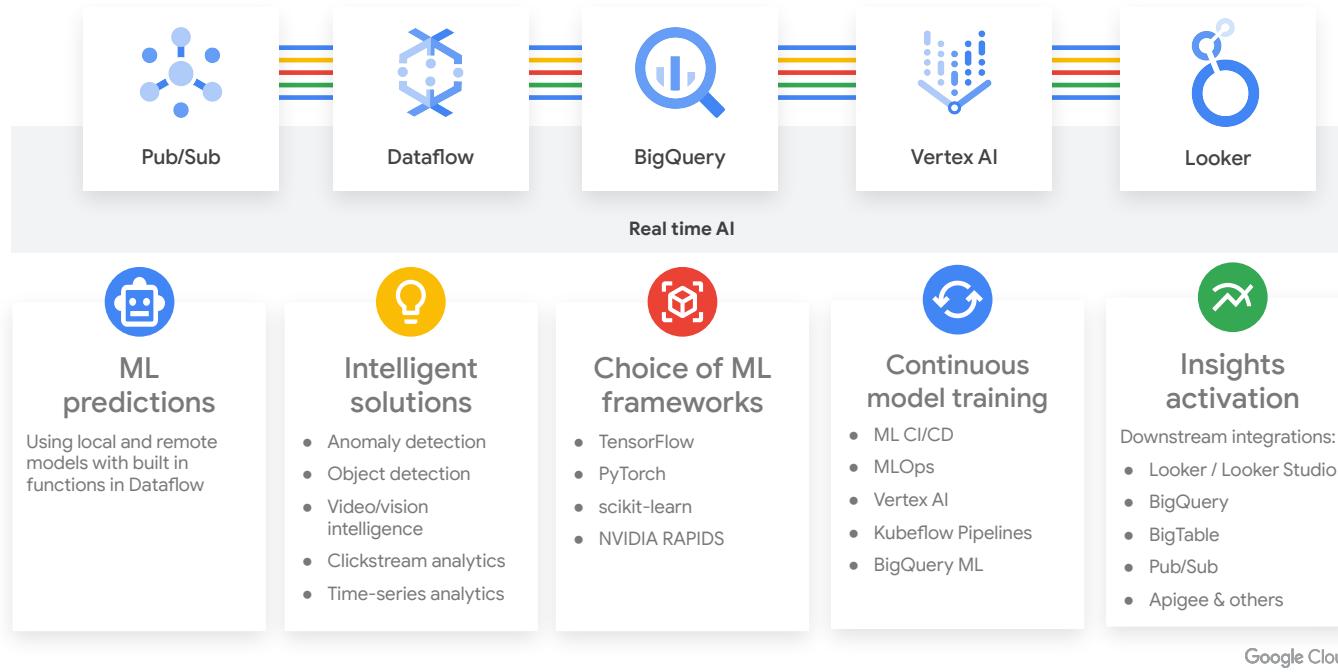
Datasets ↗		
The following limits apply to BigQuery datasets:		
Limit	Default	Note
Maximum number of datasets	Unlimited	The maximum number of datasets is limited only by storage capacity.
Number of tables per dataset	Unlimited	When creating a dataset, you can specify a maximum number of tables (50,000) or leave it at the default value of unlimited.

Google Cloud

**How can you get your data  
into BigQuery?**

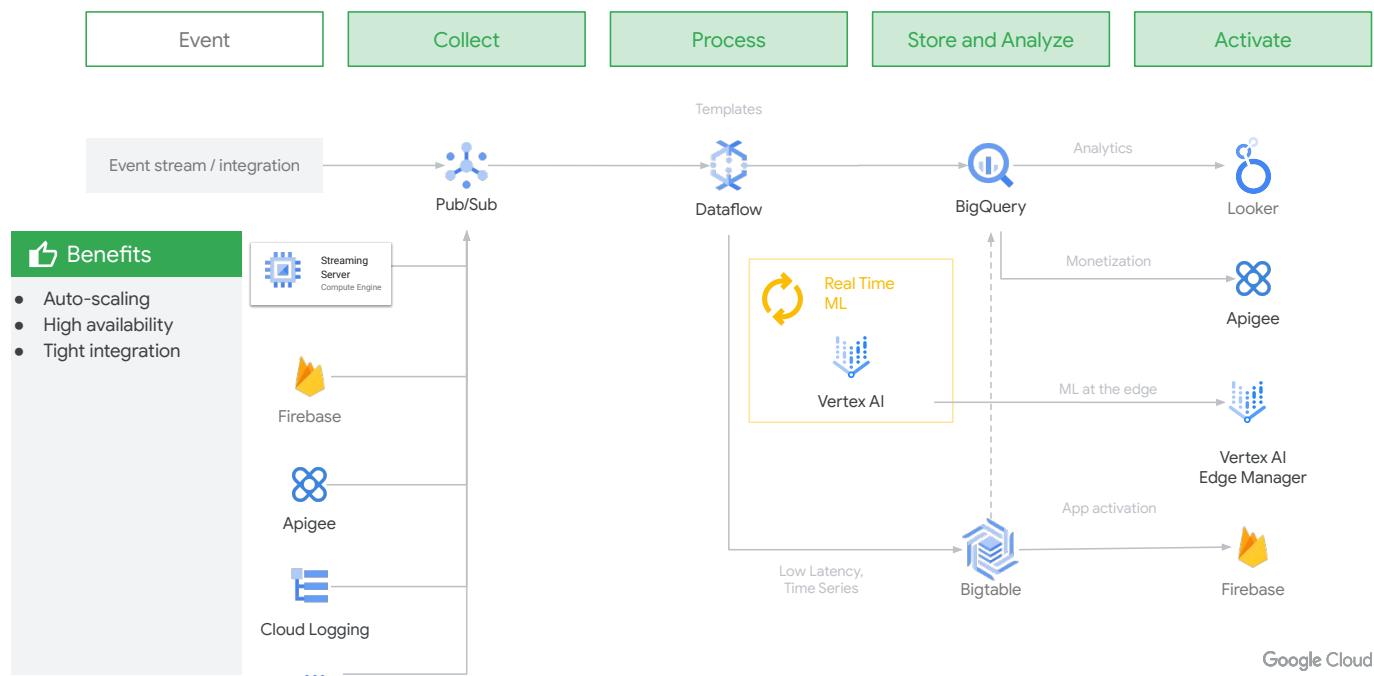
Google Cloud

## Google's Cloud Native **Golden Path** for Real Time Intelligence



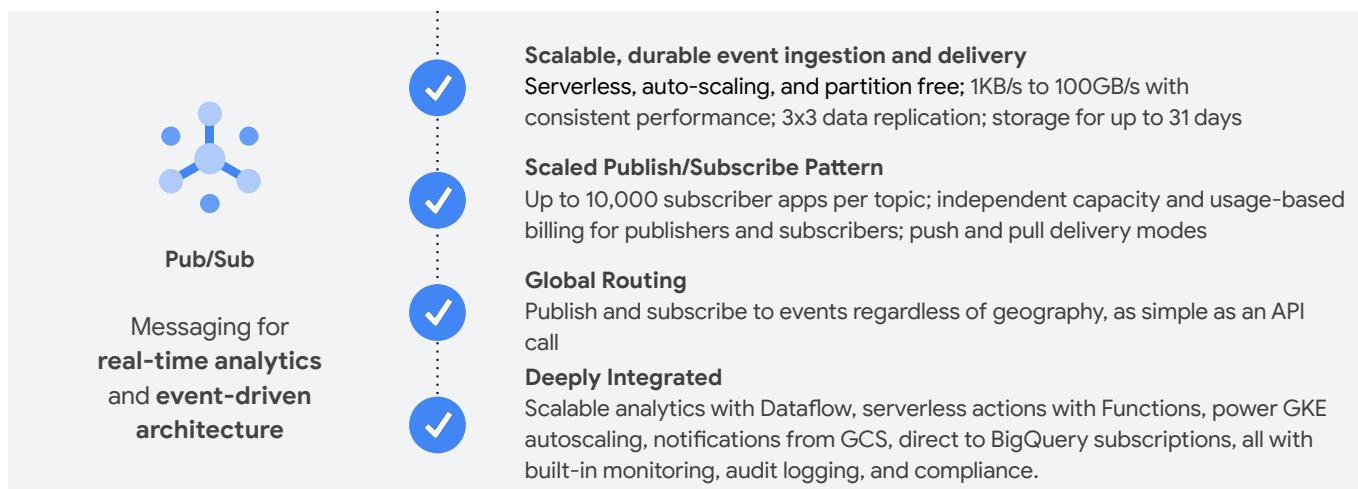
# Real Time Intelligence - High Level Architecture

Proprietary + Confidential  
Proprietary + Confidential



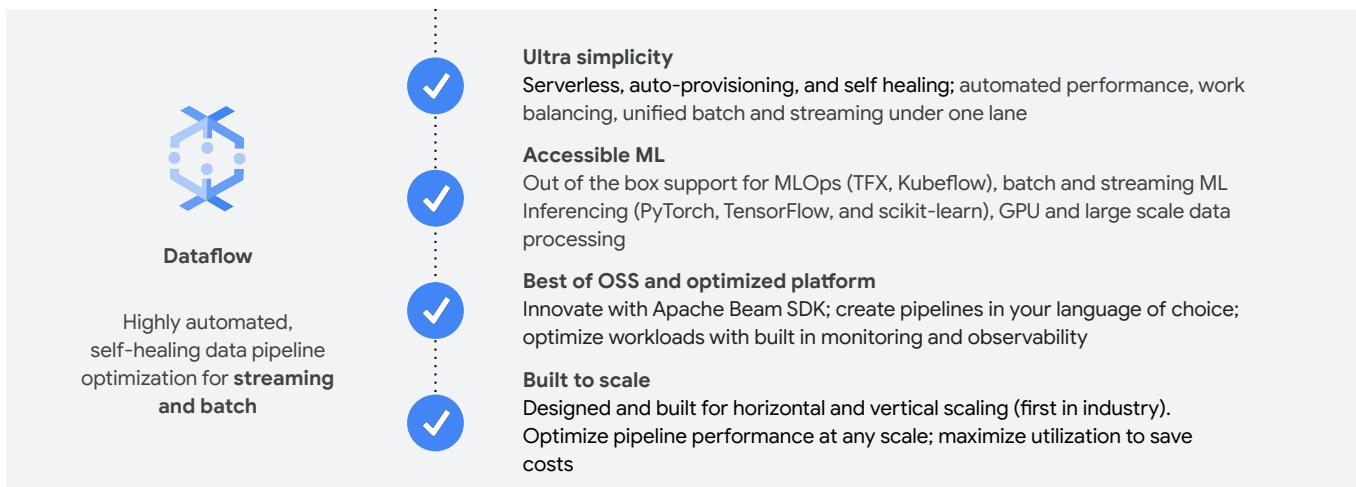
## Pub/Sub

Simple, hyperscale, durable subscription messaging when people, applications and machines need to connect insight to everything



## Dataflow

The backbone of data analytics on Google Cloud, with data pipeline optimization and automation, removing the complexity of streaming intelligence



Read more: <https://cloud.google.com/blog/topics/developers-practitioners/dataflow-backbone-data-analytics>

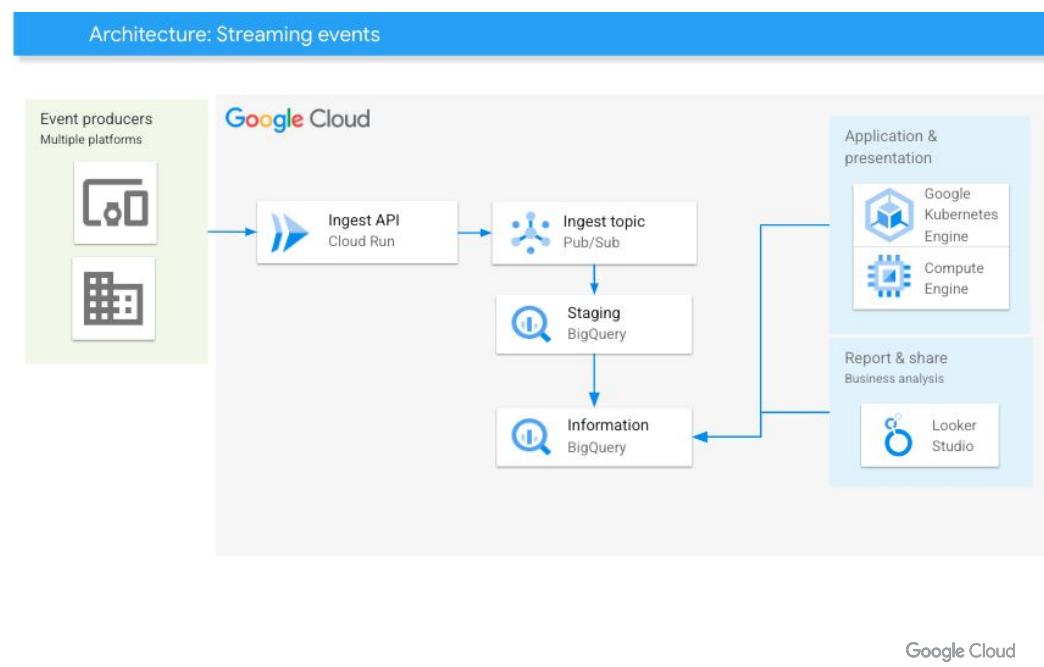
Google Cloud

## Stream events in real-time for data analytics on Google Cloud

This solution explores an example architecture pattern that ingests, processes, and analyzes a large number of events concurrently from many different sources. The processing happens as events unfold, enabling you to respond and make decisions in real-time.

Pub/Sub acts a pipeline for ingestion of real-time events. Dataflow performs the data transformation. That data is loaded into the BigQuery analytics engine.

Looker Studio or Looker can be used for creating dashboards and visualizing the data.



**Try it yourself**

**Demo :**

**[goo.gle/loff-demo24-StreamAnalytics](https://www.google.com/search?q=google+cloud+stream+analytics+demo)**

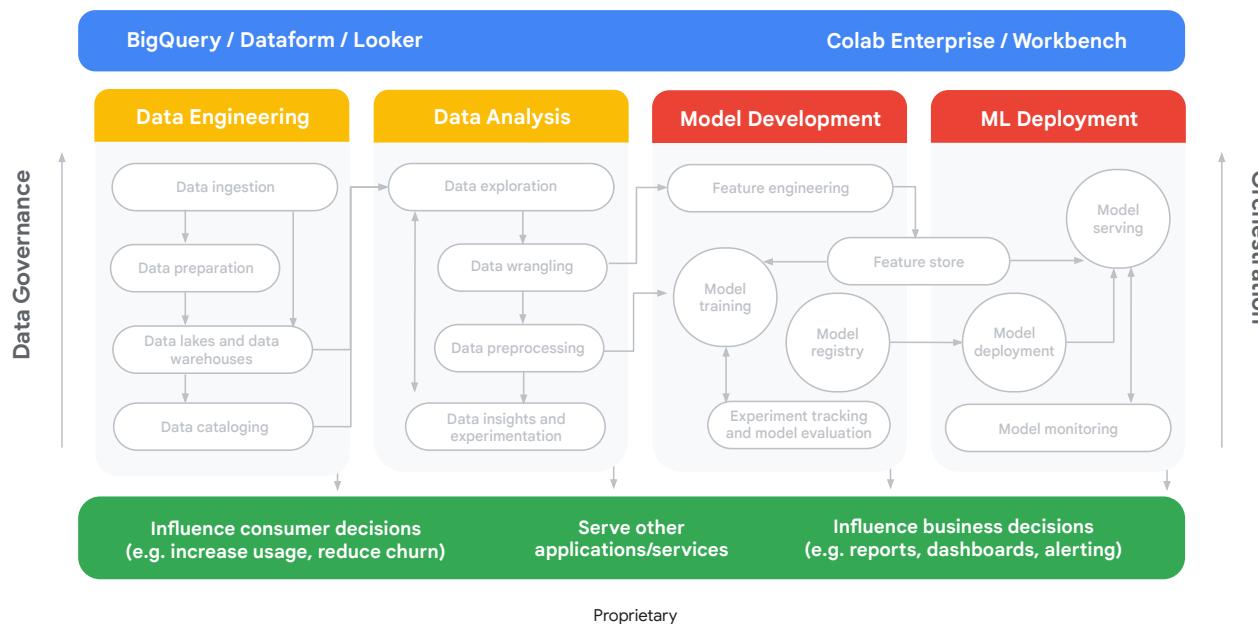
**More solutions:**

**[goo.gle/loff-demo24-click-to-deploy](https://www.google.com/search?q=google+cloud+click+deploy+solution)**

Google Cloud

# Data and AI platform

Build, experiment, deploy, and manage ML models.



**Get started  
with AI**

Google Cloud



# **What is Artificial Intelligence?**

Google Cloud

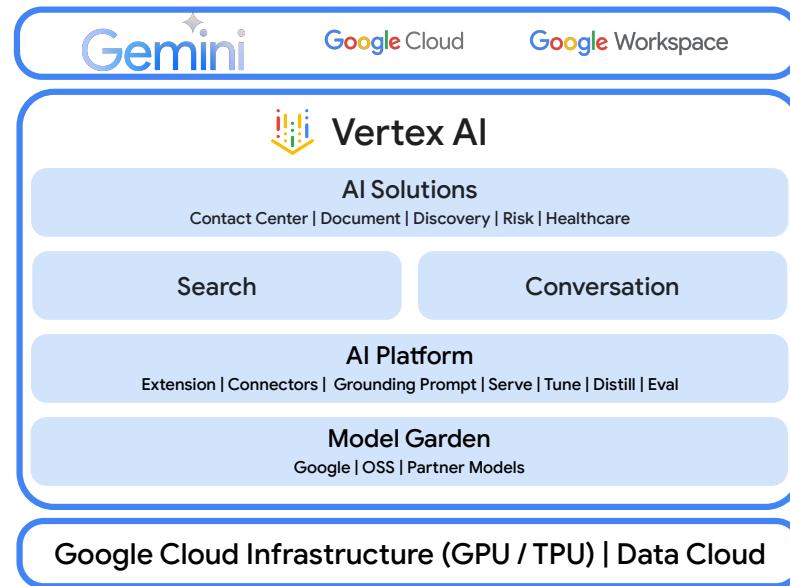
# What is Artificial Intelligence?

Artificial intelligence (AI) is a set of technologies that enable computers to perform a variety of advanced functions, including the ability to [see](#), understand and [translate spoken and written language](#), [analyze data](#), make recommendations, and more.

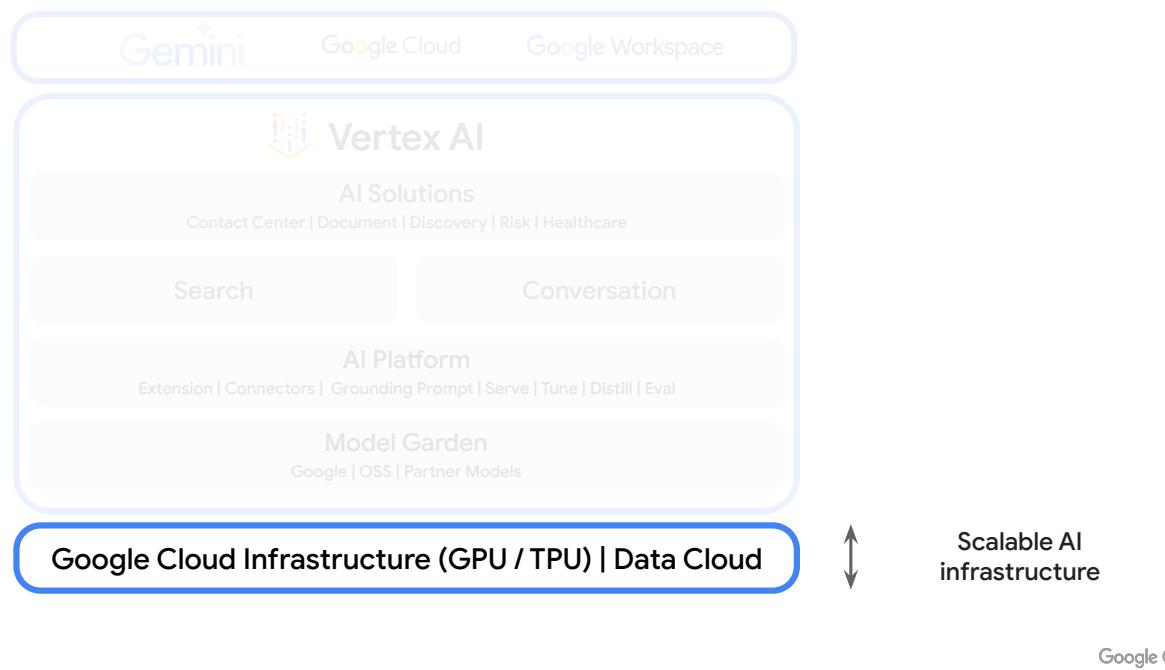
Try me!

Google Cloud

Let's take a look at how this manifests in our portfolio



Google Cloud



## Purpose-built infrastructure for Generative AI

**Tensor Processing Unit:** Designed by Google for AI at Scale



Cloud TPU v5e  
up to 2.7x inference Perf\$ vs v4

Cloud TPU v5p  
up to 2.8x LLM training vs v4

**Cost Efficient & Versatile**  
(Training & Inference)

**Powerful & Flexible**  
(full range of AI models)

**NVIDIA GPUs:** Latest NVIDIA GPUs on Google Cloud



G2 GPU VM  
2-4x performance improvement vs T4

A3 GPU VM  
3x training improvement vs A2 VM

Powered by  
NVIDIA L4 GPU

Powered by  
NVIDIA H100 GPU

- Provide a **wide variety of hardware** options
- **Scale** AI models exponentially
- Leverage our fully-managed AI platform optimized for **efficiency**
- Build with an **open source software ecosystem**

Google Cloud

# Cloud TPU v5e

Highly efficient, versatile, and scalable AI supercomputer

GA



## Cost Efficient

Up to 2.3X Training Perf/\$  
(Vs TPU v4)

Up to 2.7X Inference Perf/\$  
(Vs TPU v4)



## Versatile

Training & Inference  
(bf16, Int8, Int4)

8 VM Shapes  
(To fit diverse model sizes)

Leading AI Frameworks  
(PyTorch, JAX, TensorFlow)



## Scalable

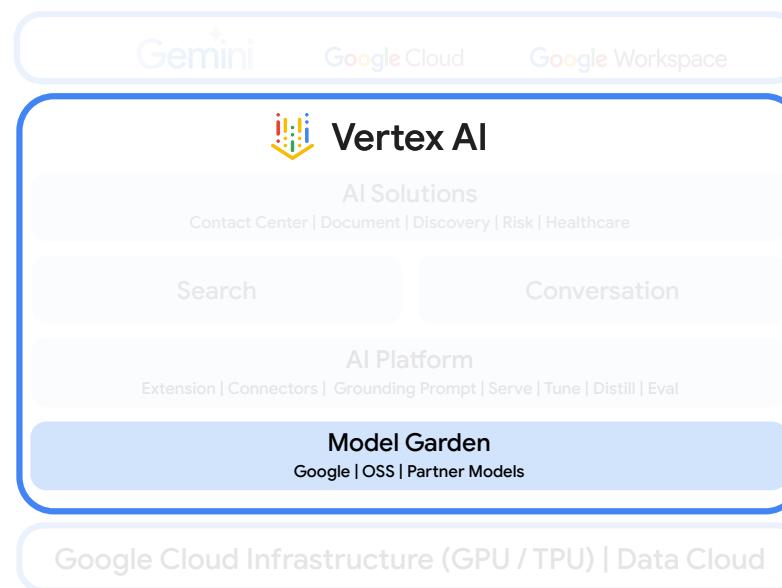
Scale to 10s Ks of Chips  
(Multislice Training)

Interchip Interconnect & DCN  
(Same XLA programming model)

Horizontal Scaling in GKE  
(Auto-provisioning & Scaling)

Training: Performance measured by Google, seq-len=2048 for GPT-3 175 billion parameter model implemented using paxml, seq-len = 2048 for 32 Billion parameter decoder model using MaxText;  
Relative performance using price of TPU v4: \$3.22/chip/hour and TPU v5e: \$1.2/chip/hour  
Inference: All results measured by Google and normalized to single-chip throughput. Precision: Llama 2 7B, 13B, 70B, GPT-J 6B: int8; GPT-J 175B, Stable Diffusion 2.1: bf16.

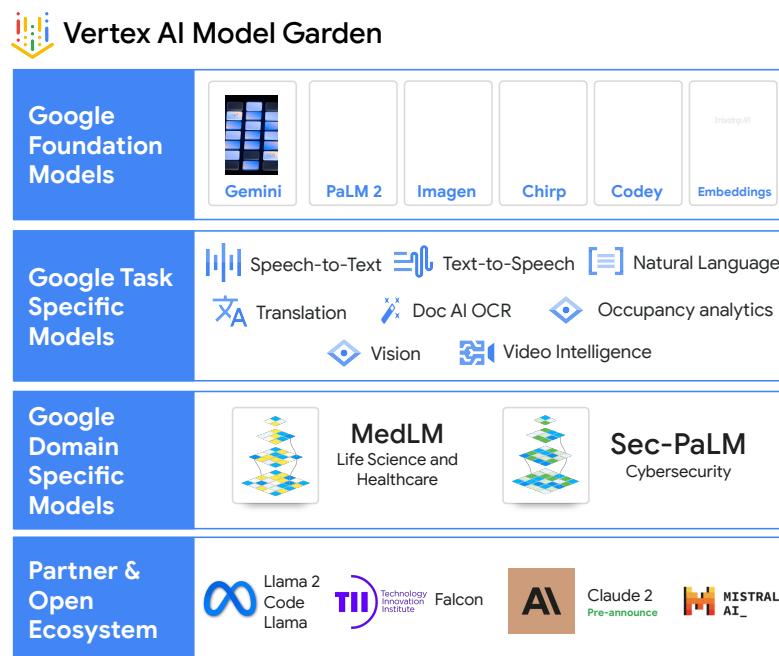
Google Cloud



↑ Access to foundation  
and open source models

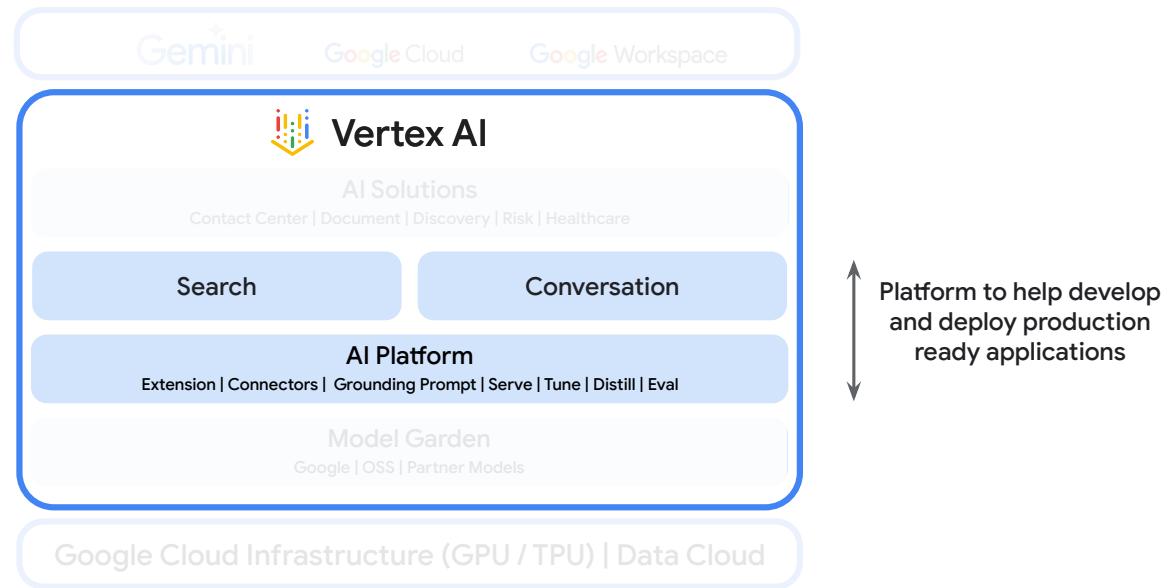
Google Cloud

## 130+ enterprise-ready foundation models in Vertex AI Model Garden



- **Choice and flexibility** with Google, open source, and third-party foundation models
- **Multiple modalities** to match every use case
- **Multiple model sizes** to match cost and efficacy needs
- **Domain-specific models** for specialized industries
- Enterprise ready with **safety, security, and responsibility**
- Decrease time to value with **fully integrated platform**

Google Cloud



## Vertex AI is built for developers



Extensive **quick start library** with code samples and jumpstarts for **developers of all levels** and ecosystems



**No cost developer labs** and training resources across Vertex products at Cloud Skills Boost



**Robust integrations** with popular third party developer tools like **Lang Chain**, **Llamaindex**, **Pinecone**, and **Weaviate**.



**Packages and extensions** to natively support Google Cloud foundation models in Google app developer frameworks like **Firebase** and **Flutter**.



Vertex AI



Colab



Flutter

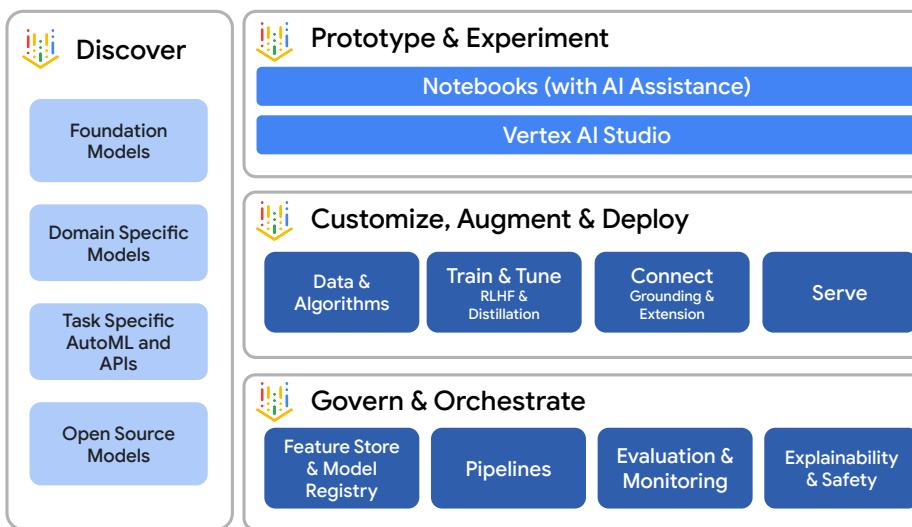


Firebase

Interfaces for  
all developers

Google Cloud

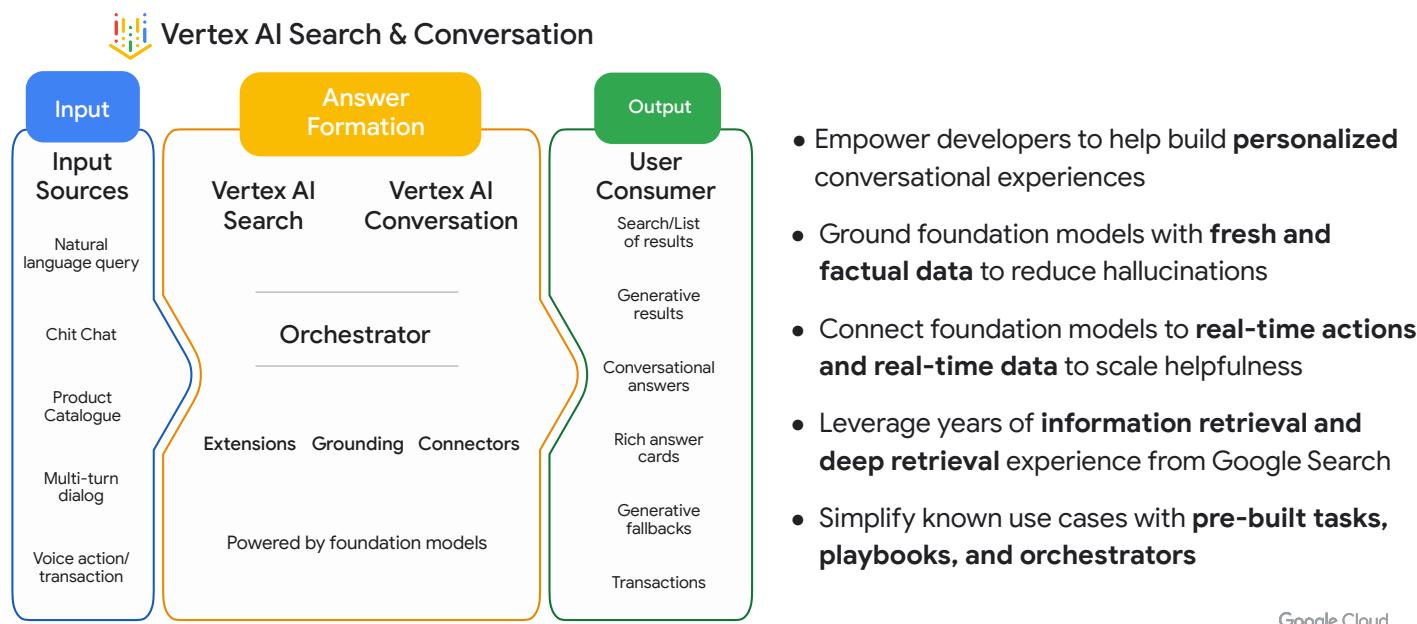
## The best of predictive and generative AI in one platform

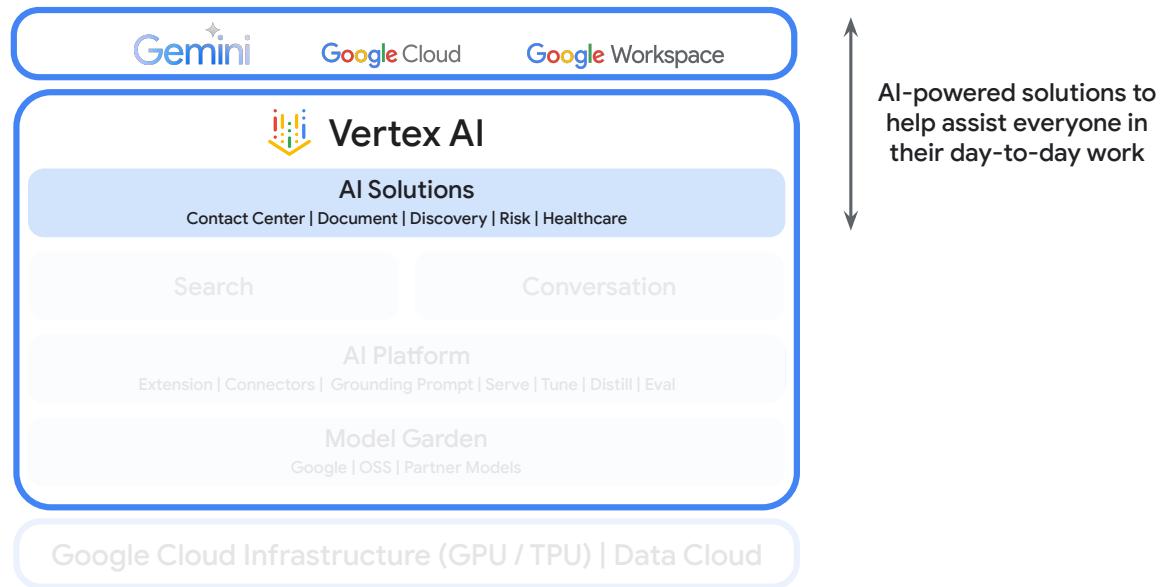


- Unified interface for **data analysis** and **AI development** to kickstart AI-powered projects across various tasks & modalities
- **Augmentation tooling** to differentiate yourself using your enterprise data
- Consistent workflows with minimal data movement **simplify data governance** for every enterprise
- Unification of **predictive and generative** through MLOps

Google Cloud

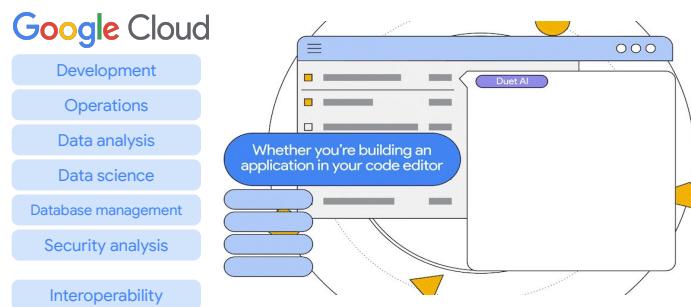
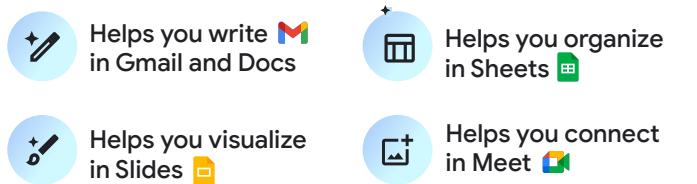
## Low code developer platform for Gen AI-powered Search and Conversation experiences





Google Cloud

## Always on AI collaborator for everyone, Gemini for Workspace



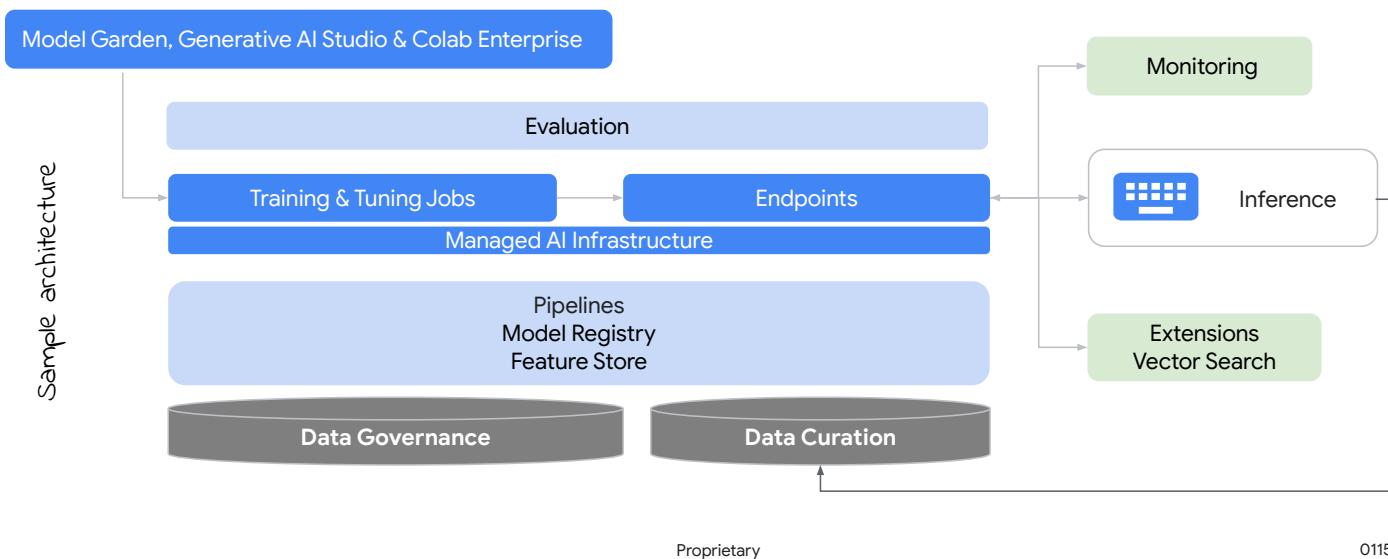
- **Democratize generative AI** for all users through **expert assistance** through Google Cloud
- Help increase employee productivity in Workspace by streamlining **creation, connection, and collaboration**
- Help drive **developer efficiency** through in the developer IDE, Cloud Console, databases, and security products
- Allow developers to focus on the **most value-add aspects** of their job

Google Cloud

# What does MLOps look like on Vertex AI?

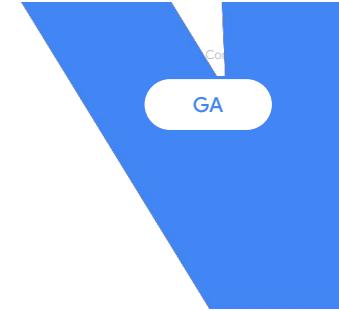
# MLOps on Vertex AI

Same Platform, Upgraded Capabilities for predictive AND generative AI



# Colab Enterprise on Vertex AI

Colab Enterprise combines the ease of use of Google Colab notebooks with the enterprise-level security and compliance capabilities of Google Cloud



Use Cases: Data science, data analysis, data engineering, ML engineering

## Collaboration & Productivity

IAM based notebook sharing  
Automatic Versioning  
Commenting (coming soon!)  
Co-editing (coming soon!)  
Generative AI powered code completion and generation

## Zero-Config & Flexible Compute

Provides both zero-config compute options, as well as access to a wide range of machine-shapes and compute

## Enterprise Ready

Will support a wide range of security and management capabilities including:

- VPC-SC
- CMEK
- Regionalization
- Cloud Monitoring
- Cloud Logging

## Available across Google Cloud

Available in BigQuery and Vertex AI (Dataproc coming soon), making it easy to work across data and AI workloads

A screenshot of a Google Colab notebook titled '3 layer neural network'. The code is written in Python and defines a function to create a three-layer neural network with specific layers and activation functions. The notebook interface includes code cells, output cells, and a sidebar with various tools and settings.



# Model Garden

Access, customize, and experiment with Google, OSS, and third party models and APIs

Model Garden [+ EXPLORE GENERATIVE AI](#) [VIEW MY MODELS](#)

Search models

Suggestions: text embedding essay outline BERT

**Foundation models**

Pre-trained multi-task models that can be further tuned or customized for specific tasks. Models marked with + are available in Generative AI Studio.

Model	Type	Description	Action
PaLM 2 for Text	+ Generative AI Language	PaLM 2 for Text	<a href="#">VIEW DETAILS</a>
Llama 2	Foundation Language	Llama 2	<a href="#">VIEW DETAILS</a>
Embeddings for Text	Foundation Language	Embeddings for Text	<a href="#">VIEW DETAILS</a>
Imagen for Image Generation	+ Generative AI Vision	Imagen for Image Generation	<a href="#">VIEW DETAILS</a>

[▼ SHOW ALL \(54\)](#)

**Fine-tunable models**

Models that data scientists can further fine-tune through a custom notebook or pipeline.

Model	Type	Description	Action
tflhub/EfficientNetV2	Classification Vision	EfficientNet V2 is a family of image classification models, which achieve better parameter efficiency and faster training speed than their predecessors.	<a href="#">VIEW DETAILS</a>
tfvision/vit	Classification Vision	The Vision Transformer (ViT) is a transformer-based architecture for image classification.	<a href="#">VIEW DETAILS</a>
tfvision/SpineNet	Detection Vision	SpineNet is an image object detection model generated using Neural Architecture Search.	<a href="#">VIEW DETAILS</a>
tfvision/YOLO	Detection Vision	YOLO algorithm is a one-stage object detection algorithm that can achieve real-time performance on a single GPU.	<a href="#">VIEW DETAILS</a>

# Deploy & Serve pre-trained and custom models

Discover and try pre-trained and open foundation AI models available on [Vertex Model Garden](#)

Carefully author an input prompt with clear, concise and informative instructions on [Vertex Generative AI Studio](#)

Deploy to fully managed endpoints with [Vertex Predictions](#) supporting autoscaling, private endpoints, and wide selection of CPUs & GPUs

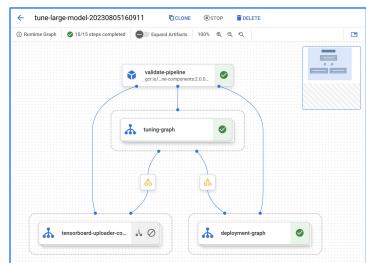
- Model co-hosting
- Custom prediction routines



Powered by AI Infrastructure for your generative AI needs  
NVIDIA GPUs A100, L4 and H100 | TPU v3, v4 and v5e

# Managing artifacts

Orchestrate & manage training & tuning jobs via [Vertex AI Pipelines](#) to trace lineage from dataset to model



Manage predictive & generative AI models in [Vertex AI Model Registry](#) with access to evaluation metrics and one-click deployment

Type	Source
Imported	Model Garden
Imported	Custom training
Tabular	AutoML training
Large model	Generative AI Studio

Proprietary

Store, manage, and serve features including embeddings in [Vertex AI Feature Store](#) to support predictive & generative AI needs

Feature Table			
CREATE			
All	307	Search	
Feature Table		Name	Feature type
User	17	FeatureTable1	--
Destination	8	text_snippet	string 12345
Fare	11	text_snippet_e	Embedding 12345

0119

# Evaluation services for predictive AI



## Evaluate

Iteratively and independently run model evaluations on new datasets at scale



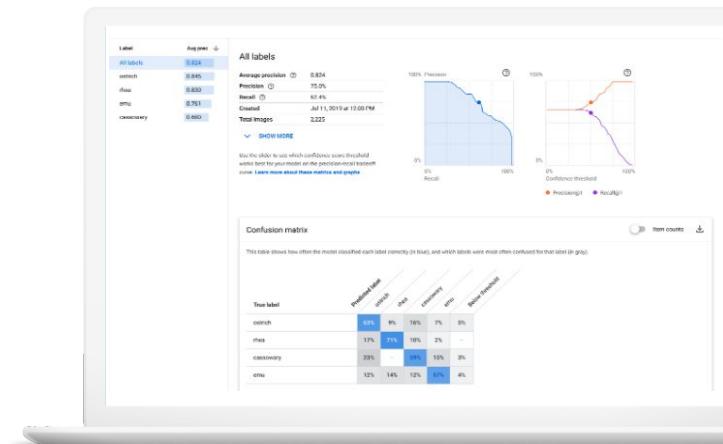
## Compare

Visualize and compare different model evaluations to identify the best model for production deployment



## Assess

Assess the performance of models on different slices and evaluated annotations



# Evaluation services for generative AI

1

## Automatic Metrics

rougeLSum	0.1345
bleu	0.763

Assess the performance of a model with task-specific metrics computed based on reference data

- Fast and efficient
- Standard method used in academia and many open benchmarks

2

## AutoSxS

model_a_win_rate	0.35
model_b_win_rate	🏆 0.65

Compare the performance of 2 models with an arbiter model

- Consistent performance with human evaluation, while being faster, cheaper, available on demand

3

## Safety Bias

Safety attribute: Toxicity	
GenderID: female	0.89
GenderID: male	0.12
GenderID: transgender	0.02
GenderID: non-binary	0.67

Understand if the safety profile of your model is biased against a certain identity group.

# Model Monitoring

## Monitor & Explain

- Model monitoring for skew & drift detection
- Feature-based explainability
- Example-based explainability

## Safety Scores

Score “harmful categories” and topics that may be considered sensitive.

```
{
  "predictions": [
    {
      "safetyAttributes": {
        "categories": [
          "Hate",
          "Toxic",
          "Violent",
          "Sexual",
          "Insult",
          "Obscene",
          "Death, Harm & Tragedy",
          "Firearms & Weapons",
          "Public Safety",
          "Health",
          "Religion & Belief",
          "Drugs",
          "War & Conflict",
          "Politics",
          "Finance",
          "Legal"
        ],
        "scores": [
          0.7,
          0.4,
          0.6,
          0.8,
          0.6,
          0.6,
          0.6,
          0.6,
          0.6,
          0.6,
          0.6,
          0.6,
          0.6,
          0.6,
          0.6,
          0.6
        ]
      }
    }
  ],
  "content": "<>"
}
```

## Recitation Checks

Help ensure that our models do not replicate existing content at length



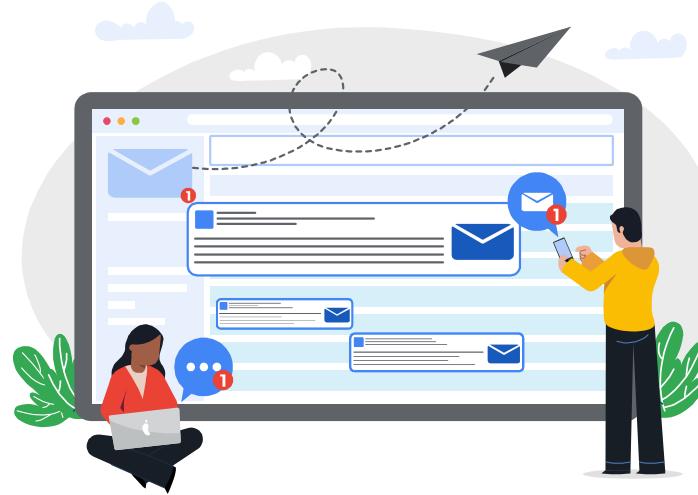
# Security with Google Cloud AI

## Data governance during generative AI model usage

Google processes prompts to provide the service.

Prompts input to the foundation model to generate a response, are encrypted in transit.

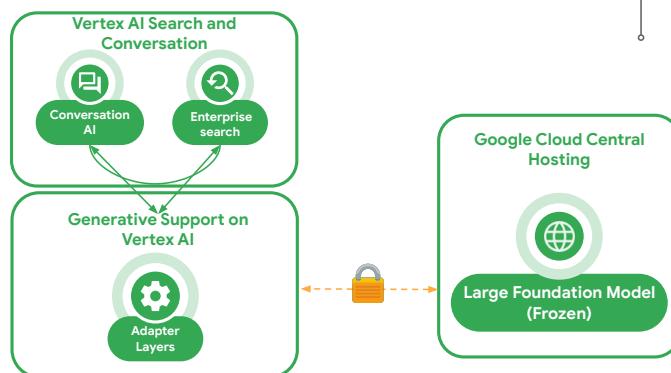
Google does not use prompt data to train its models without the express consent of its customers.



Google Cloud

## Security during Model Tuning

- Parameter Efficient Fine Tuning (PEFT)
- Customer specific adapter weights
- Foundational model remains frozen during inference



- Input data is stored securely
- Adapter weights are stored securely
- Customer can delete adapter weights at any time
- Customer Data will not be logged to train foundation models by default



Google Cloud

## Security Controls for Generative AI



**VPC Security Controls:** Help protect against accidental or targeted action by external entities or insider entities, to minimize unwarranted data exfiltration risks.



**Customer Managed Encryption Keys:** Control the keys that protect your data at rest. Support for External Key Management



**Access Transparency:** Near real-time logs offer insight when Google Cloud administrators access your content - including access justification.





## Technical safeguard: Safety filters

**Safety attributes** label “harmful categories” and topics that may be considered sensitive.

Each safety attribute has an associated confidence score in [0.0, 1.0], rounded to one decimal place, reflecting the likelihood of the input or response belonging to a given category.

You can test Google's safety filters and define confidence thresholds that are right for your business and take comprehensive measures to detect content that violates Google's usage policies or terms of service.

Google Cloud



## Recitation checks

Recitation checks help ensure that our models do not replicate existing content at length

We've designed our systems to limit the chances of replicating existing content at length and we will continue to improve how these systems function.

Generally, if our API does directly quote at length from a page on the web that was part of our pre-training dataset, it cites that page in our output from the model, or if over a certain length, it will block the output.



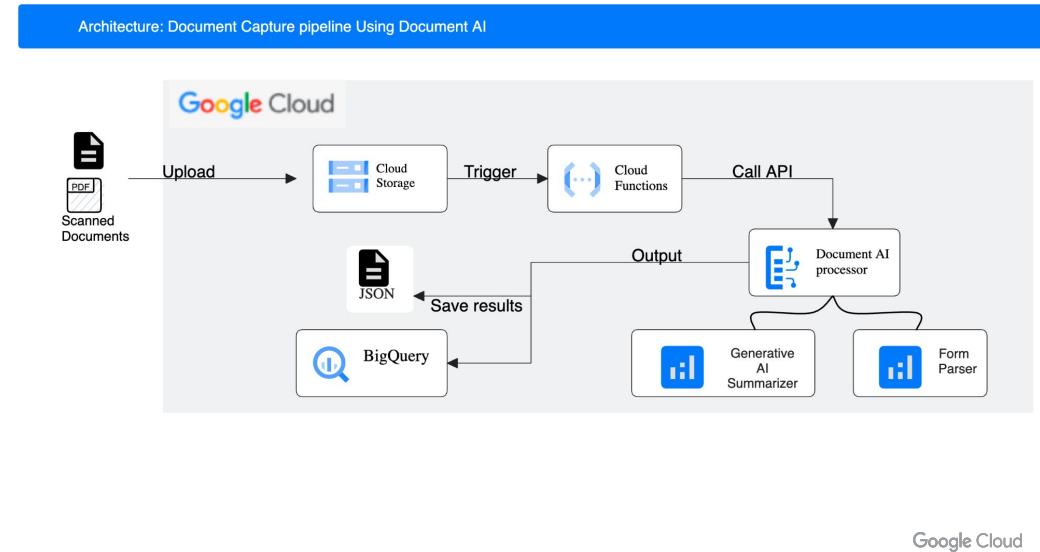
Google Cloud

## Extract data from your documents using Generative AI

Automatically extract structured data (e.g., from forms, invoices) using Document AI's form processor.

The solution leverages genAI to produce concise summaries of the processed documents, highlighting key information for enhanced understanding.

Extracted data and summaries are stored in BigQuery, enabling powerful analysis, visualization, and machine learning applications.



## Try it yourself

Demo :

[goo.gle/loff-demo24-DocAIForm](https://goo.gle/loff-demo24-DocAIForm)

More solutions:

[goo.gle/loff-demo24-click-to-deploy](https://goo.gle/loff-demo24-click-to-deploy)

Google Cloud

# Try it yourself

## [Learn more about Generative AI](#)

The screenshot shows the Google Cloud Skills Boost interface. At the top, there are navigation links: Google Cloud, Paths (which is highlighted in blue), Explore, and Subscriptions. Below this, the title "Google Cloud Skills Boost" is displayed. The main content area features a large card titled "Introduction to Generative AI Learning Path". It includes a small icon of a brain, the number "5 activities", the text "Last updated 4 months", and "Managed by Google Cloud". A descriptive text states: "This learning path provides an overview of generative AI concepts, from the fundamentals of large language models to responsible AI principles." Below this text is a "Start learning path" button. At the bottom of the card, there are three smaller cards with titles: "01 Introduction to Generative AI", "02 Introduction to Large Language Models", and "03 Introduction to Responsible AI".

## [Train a model using Vertex AI and Python SDK](#)

The screenshot shows a Jupyter Notebook interface. The top menu bar includes File, Edit, View, Run, Kernel, Git, Tabs, Settings, and Help. On the left, there's a file browser with a search bar and a sidebar with icons for Notebooks, Console, and another Python 3 kernel. The main area displays a list of files in the current directory: "src" and "tutorials". Both files were last modified 18 hours ago. A red box highlights the "Python 3" icon in the sidebar.

Google Cloud

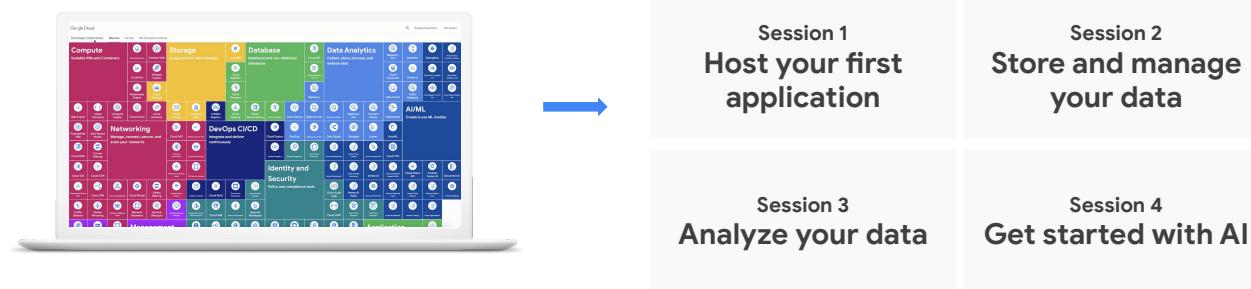
**What's next?**

Google Cloud



# What is Google Cloud?

It's how you can...





# Thank you

Check out more ways to learn  
in the description below

Google

Proprietary + Confidential

# Graveyard

Google

Limitless data

## Limitless data scale

275 EB

analyzed across BigQuery  
in December 2021

110+ TB

data analyzed per second

“

...BigQuery continued to **scale** in storage, compute, concurrence, ingest and reliability as we added **more and more** users, traffic, and data.”

Nikhil Mishra  
Sr. Director of Engineering, Verizon Media

verizon<sup>✓</sup>



Google Cloud

## Host a Highly Available SQL database on Google Cloud

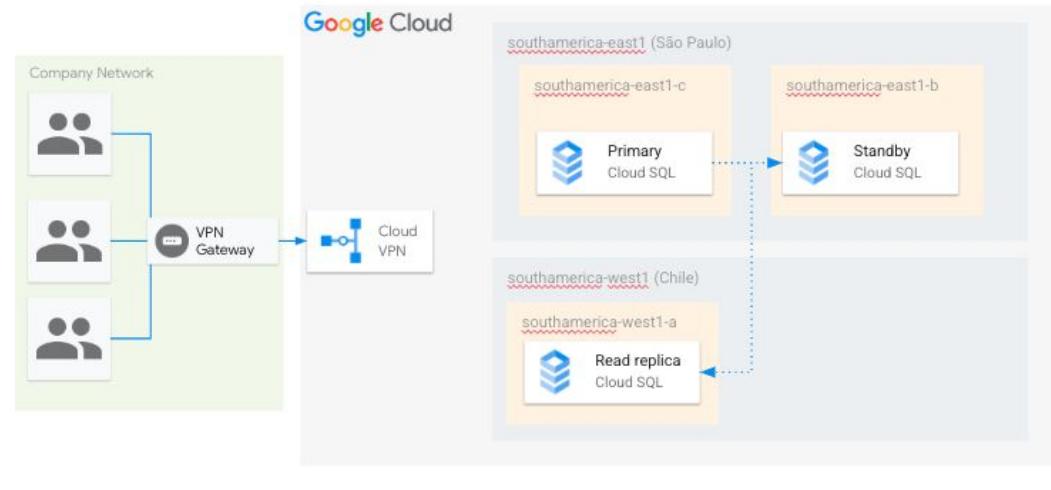
Cloud SQL instance with high-availability and cross-region replica.

Leverages regional and zonal failovers as well as read replicas offered by Cloud SQL

Uses Cloud VPN to connect company network to Google Cloud

Suitable for a wide range of applications and industries, including e-commerce platforms, financial systems, customer relationship management (CRM) tools, and more.

Architecture: Cloud SQL with HA and DR



Google Cloud

**Try it yourself**

**Demo :**

**[goo.gle/loff-demo24-CloudSQLMR](https://www.google.com/search?q=google+CloudSQL+ML+Demo)**

**More solutions:**

**[goo.gle/loff-demo24-click-to-deploy](https://www.google.com/search?q=click+to+deploy+CloudSQL+ML+Demo)**

Google Cloud