# Introduction to Distributed workload with Ray on Kubernetes

Mofi Rahman

@mofi.dev

Abdelfettah Sghiouar

@boredabdel.bsky.social

# Distributed Computing

Why use one computer to solve a problem when you can use thousands?

# Distributed Computing

**Python** is the "lingua franca" of AI

With GenAI distributed compute is no longer optional, it is **required**

# Why Distributed Computing?

Scalability

Availability

Efficiency

Flexibility

# Challenges

Consistency

Fault Tolerance

Concurrency Control

Load Balancing

Security Concerns
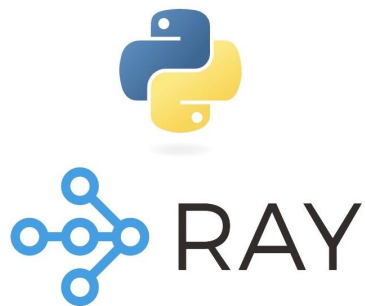
Complexity of Management

# CAP Theorem

Consistency
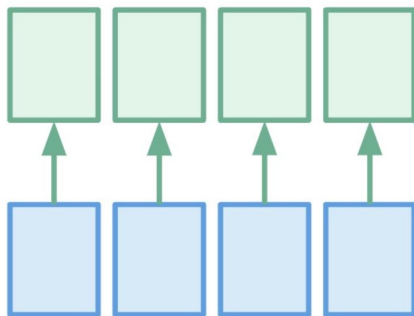
Availability

Partition Tolerance
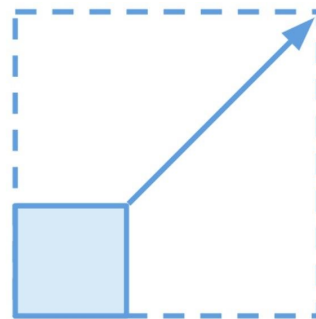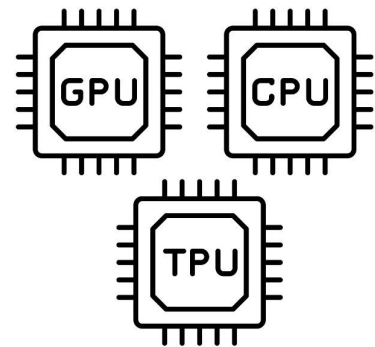
# Ray

# Ray - Key Characteristics

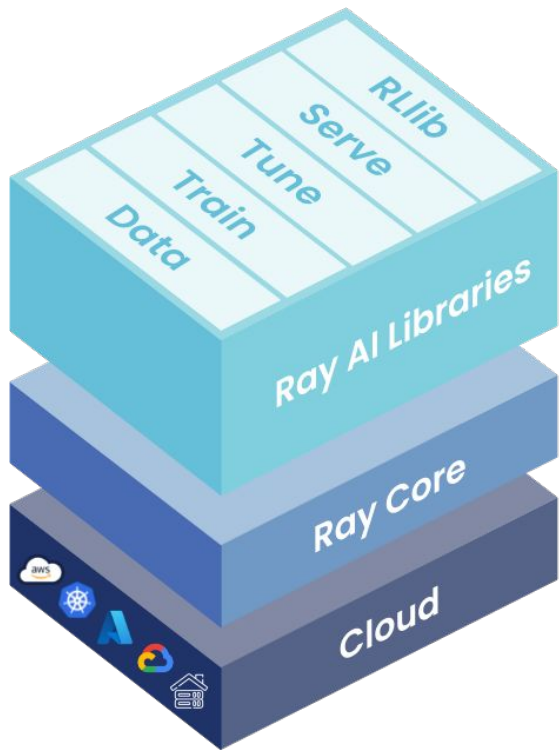Python first approach, open source

Simple and flexible API

Scalability

Support for bleeding edge hardware

GPU CPU TPU

Google Cloud

# Ray - Components



high-level libraries that enable simple scaling of AI workloads

a low-level distributed computing framework with a concise core and Python-first API

Google Cloud

# Ray AI Libraries
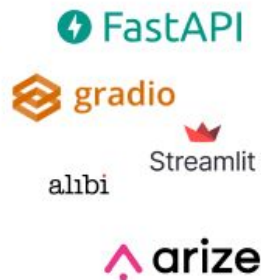


Ray Data · Ray Train · Ray Tune · Ray Serve

Data Preprocessing → Training → Tuning → Serving

High-level libraries that make scaling easy for both data scientists and ML engineers.
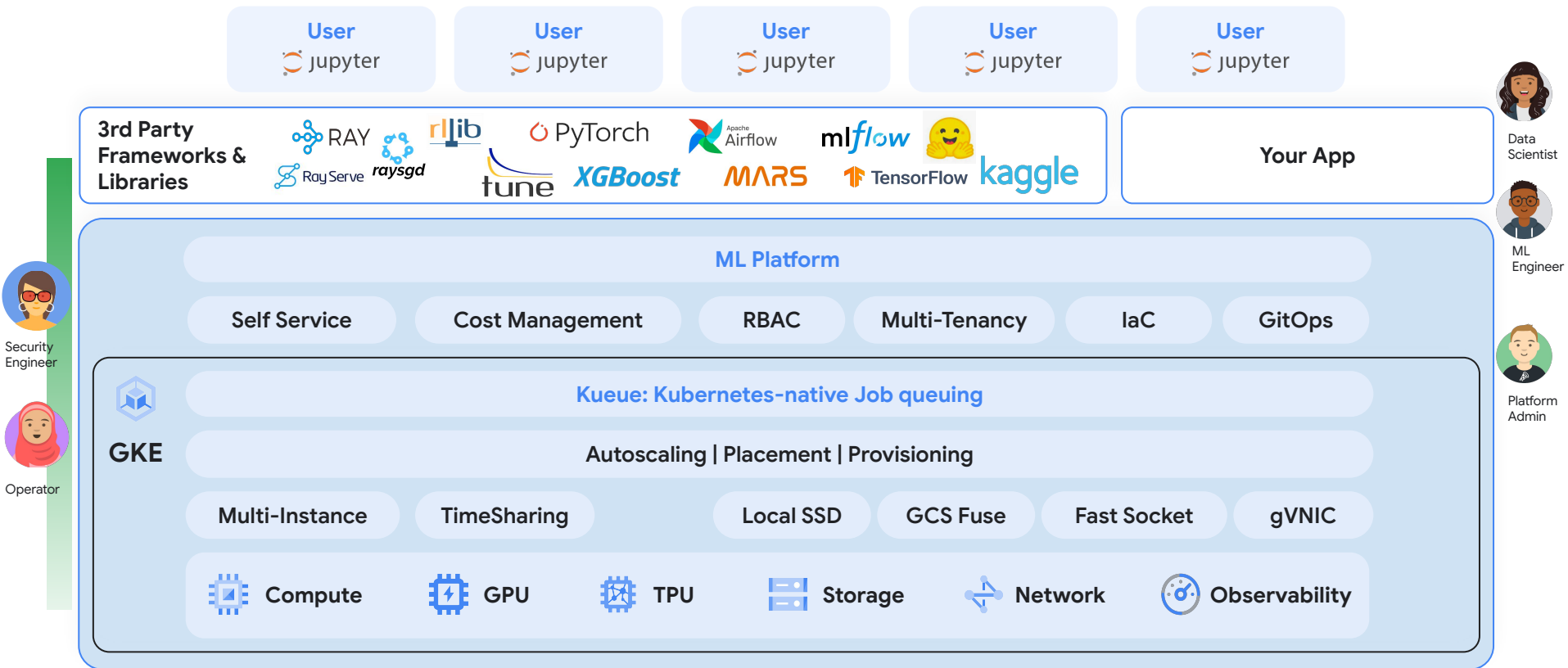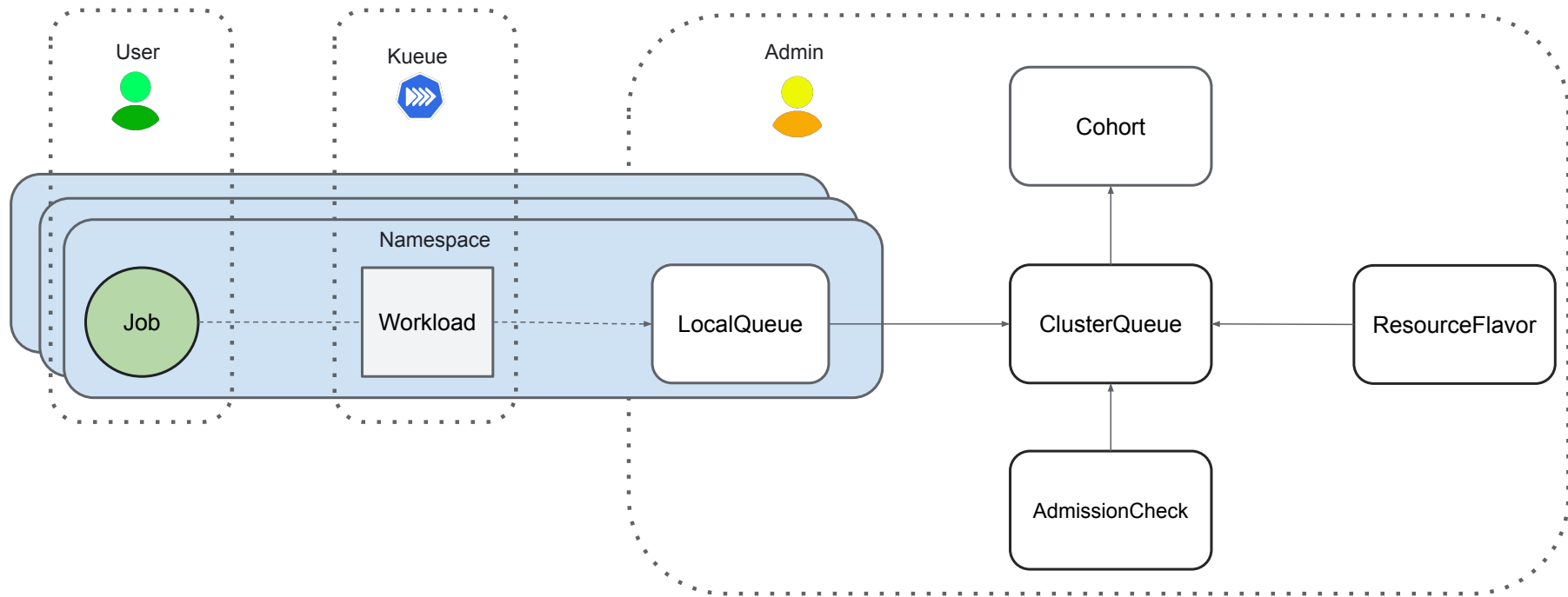
# Key Concepts

Tasks

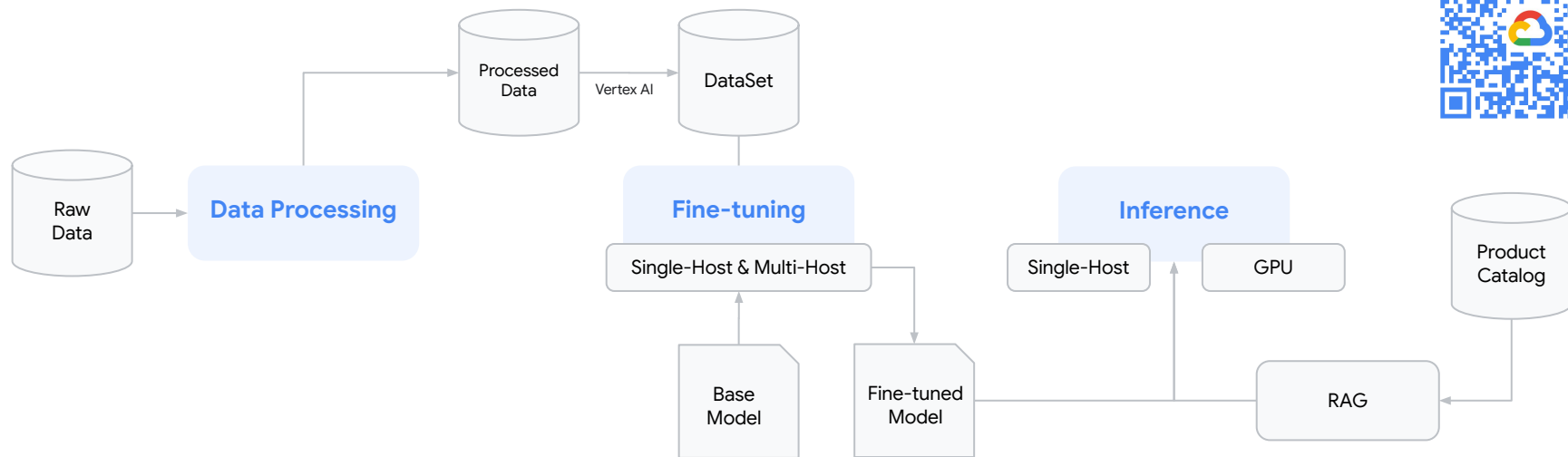Actors

Objects

# Demo

# ML Platform

# The ever growing AI/ML ecosystem

**User** Jupyter
**User** Jupyter
**User** Jupyter
**User** Jupyter
**User** Jupyter

**3rd Party Frameworks & Libraries**

RAY
Ray Serve
raysgd
rllib
tune
PyTorch
XGBoost
Apache Airflow
MARS
mlflow
TensorFlow
kaggle

**Your App**

Data Scientist

**ML Platform**

| Self Service | Cost Management | RBAC | Multi-Tenancy | IaC | GitOps |

ML Engineer

**GKE**

**Kueue: Kubernetes-native Job queuing**

**Autoscaling | Placement | Provisioning**

| Multi-Instance | TimeSharing | Local SSD | GCS Fuse | Fast Socket | gVNIC |

Platform Admin

Compute · GPU · TPU · Storage · Network · Observability

Security Engineer

Operator

Google Cloud

# Kueue

User

Kueue

Admin

Namespace

Job

Workload

LocalQueue

Cohort

ClusterQueue

ResourceFlavor

AdmissionCheck

Google Cloud

17

# Data Preprocessing, LLM Fine-Tuning, Inference at Scale

# Data Preprocessing, LLM Fine-Tuning, Inference at Scale

# So how do we run Ray?

# Ray can be run anywhere!

| Ray AI Libraries | | | |
|---|---|---|---|
| Data | Train | Tune | Serve |

| Ray Core |
|---|

| Containers | VM Launcher |
|---|---|
| Kubernetes | Cloud VMs |

# Enter Kuberay

| Ray AI Libraries | | | |
|---|---|---|---|
| Data | Train | Tune | Serve |

| Ray Core |
|---|

| KubeRay | VM Launcher |
|---|---|
| Kubernetes | Cloud VMs |

Manages the lifecycle of Ray clusters and associated applications on Kubernetes.

# Why Ray on Kubernetes?

Infrastructure Automation

Scalability

High Availability and Reliability

Advanced Device Management

Complex, multi-cloud deployments

# KubeRay: The best solution for Ray on Kubernetes

Ray

Create / Update

Health check

Scaling requests for Tasks / Actors

Monitoring

KubeRay

Read

Create / Delete

Update (observability)

Kubernetes

# KubeRay APIs

## RayCluster

Manage and scale Ray clusters

Ideal for prototyping / development

## RayJob

Execute a Ray job with ephemeral Ray clusters

Ideal for productionizing Ray batch workloads

## RayService

Deploy a Ray Serve application with zero-downtime upgrades

Ideal for inference in production

# Demo

# **Ephemeral vs Persistent Ray Cluster**

**Prons**

- Reproduciblity
- No need for maintenance since a fresh cluster is started for each job
- Better observability for single distributed jobs metrics

**Cons**

- Startup latency can be long
- Ray dashboard lasts as long as the workload
- Logs and Metrics has to be stored outside the cluster

# Ephemeral vs Persistent Ray Cluster

**Prons**

- Startup latency for workload is small
- Minimal packaging required if the clusters already has the dependencies
- Ray dashboard can be used to track history

**Cons**

- New dependencies are tricky
- If cluster is brought down for maintenance the behavior can be unpredictable

# Security

- Ray endpoints are not locked down by default (not Authn or Authz)
  - ⇨ Leverage Cloud Providers tools to secure the endpoints (LBs, proxies...)

- A lot of layers to secure (Kubernetes, Ray pods, Workloads...)
  - ⇨ Kubernetes tools can be used to secure the clusters and the pods (RBAC, namespaces, Quotas, Pod Security Policies...)

# KubeRay is growing!

## Community

1000+ commits

140+ contributors

## Adoption

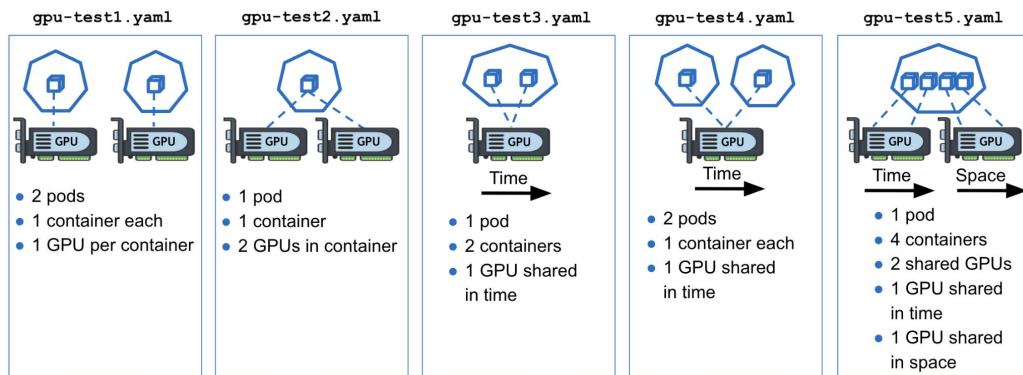100s organizations

50+ blogs & talks

## Scale

10K Ray clusters

40K Pods

# DRA: Optimizing Resource Allocation

DRA enhances the Kubernetes scheduler with awareness of Ray's needs and the dynamic nature of certain workloads:

- Optimized resource utilization
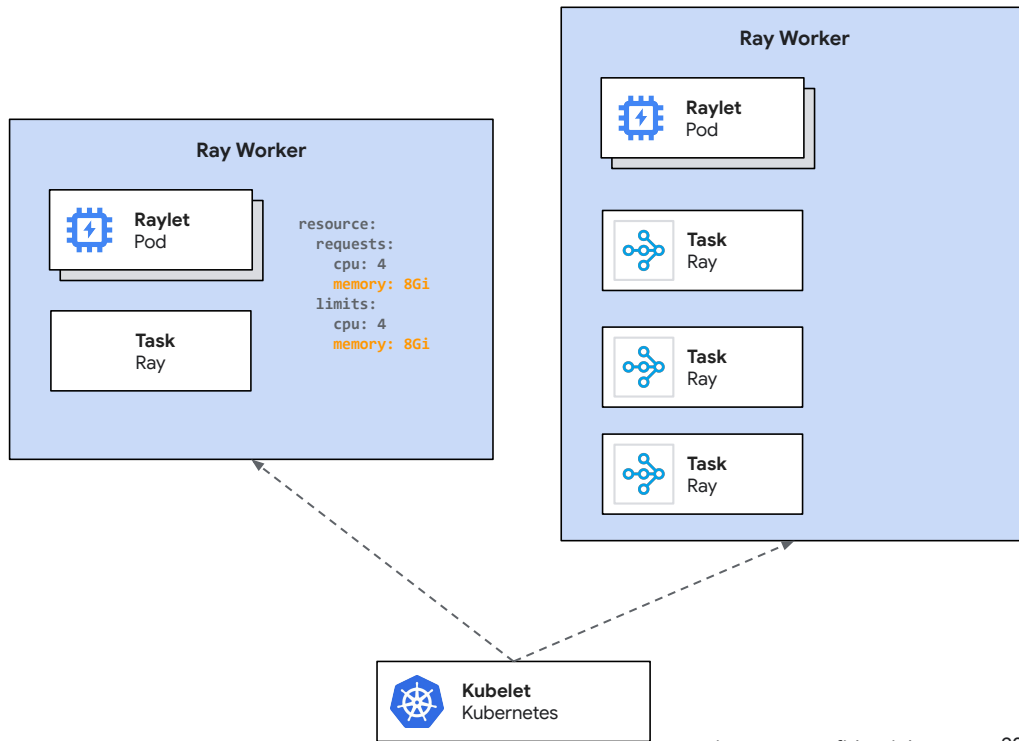- Improved cluster efficiency



**Kubernetes v1.31 introduced new Device Resource Assignment (DRA) APIs.**

# In-place VPA: Minimizing Disruptions

In-place Vertical Pod Autoscaling enables elastic memory consumption for Ray containers without requiring restarts.

Prevent performance degradation and risk of OOM-kill with resizable Pod memory
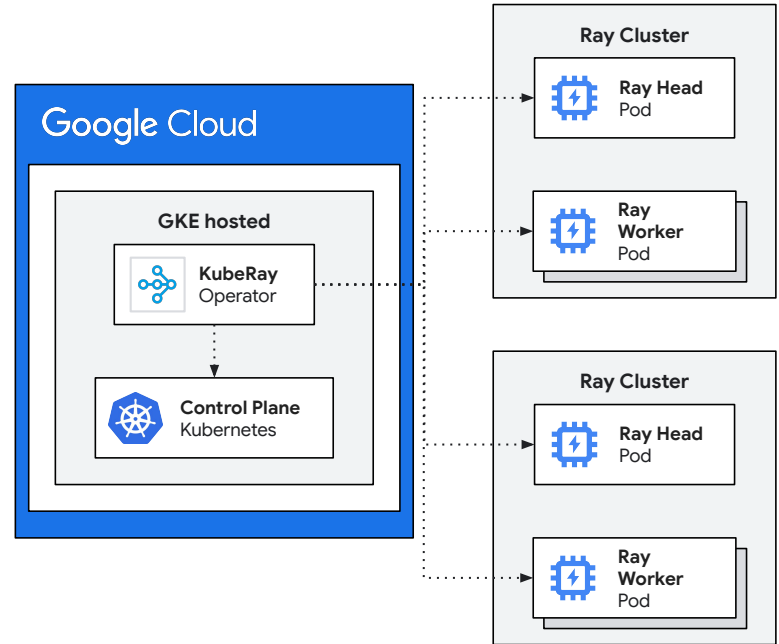
# Ray Operator on GKE

GKE hosts the KubeRay operator on your behalf.

Optimally configured for high performance and scalability.

Get started with Ray on GKE with a single option using:

- gcloud CLI
- Google Cloud Console
- Terraform

See Enable the Ray operator on Google Kubernetes Engine for more details.

# Feedback 🙏

# Thank you

Google Cloud