

# Comprehensive Implementation Guide: PRJ-GCP-AI-082 - Custom Model Training with Vertex AI

---

**Author:** Manus AI **Date:** January 26, 2026

---

## 1. Project Overview

---

The PRJ-GCP-AI-082 project establishes a **secure, compliant, and production-ready MLOps framework** on Google Cloud Platform (GCP) utilizing the comprehensive capabilities of **Vertex AI**. The core objective is to facilitate custom machine learning model training and deployment within a hardened enterprise environment, ensuring strict adherence to responsible AI principles, stringent data privacy regulations, and robust security standards.

This solution is designed for organizations operating in regulated industries or those deploying high-impact AI systems where governance and auditability are paramount. It integrates key Vertex AI services—including Pipelines, Training, Model Registry, and Monitoring—to create an end-to-end governed workflow. The architecture is specifically engineered to mitigate common risks such as data exfiltration and model bias, providing a foundation for trustworthy and scalable AI operations.

Component	Purpose	Key Feature
Cloud Platform	Google Cloud Platform (GCP)	Enterprise-grade security and global scale.
Core Services	Vertex AI (Pipelines, Training, Monitoring)	Unified platform for MLOps lifecycle management.
Data Security	Cloud Storage, BigQuery, VPC Service Controls	Secure, encrypted data hosting and perimeter defense.
Governance	Vertex AI Model Cards, Explainable AI	Transparency, auditability, and responsible AI enforcement.
Tools	gcloud CLI, gsutil CLI, Docker	Command-line and containerization tools for deployment.

## 2. Business Context

---

The proliferation of AI models in enterprise operations introduces significant challenges related to **data privacy, model security, and regulatory compliance**. Without a structured MLOps framework, organizations face substantial compliance and ethical risks, particularly when models influence critical business decisions. This project directly addresses these challenges by embedding governance and security into the core MLOps workflow.

### The Problem and Solution

**The Problem:** Many organizations struggle to leverage the full potential of GCP AI services due to a lack of a secure, compliant MLOps infrastructure. This gap leads to:

- 1. Compliance Risk:** Difficulty in proving model fairness, lineage, and adherence to regulations like GDPR or the EU AI Act.
- 2. Security Risk:** Potential for data leakage and unauthorized access to sensitive training data and model artifacts.
- 3. Operational Inefficiency:** Manual, non-repeatable processes for model deployment that lack audit trails.

**The Solution:** PRJ-GCP-AI-082 implements a secure MLOps pipeline using **Vertex AI with responsible AI controls, continuous model monitoring, and data protection.** The solution enforces AI governance, includes mechanisms for bias detection, and provides comprehensive logging for complete auditability, ensuring the entire machine learning lifecycle is transparent and compliant from inception to production.

## Quantified Business Value and ROI

The implementation of this secure MLOps framework translates directly into quantifiable business value across several dimensions:

Value Proposition	Description	Quantified Impact / ROI
<b>Risk Mitigation &amp; Compliance</b>	Proactive adherence to global AI regulations (e.g., EU AI Act, NIST AI RMF) and data privacy laws (GDPR, CCPA).	<b>Up to 90% reduction</b> in compliance-related fines and legal costs. <b>100% audit readiness</b> for high-risk AI systems.
<b>Operational Efficiency</b>	Automated, repeatable ML pipelines replace manual deployment processes.	<b>30-50% faster time-to-market</b> for new models. <b>20% reduction</b> in MLOps engineering overhead.
<b>Responsible AI &amp; Trust</b>	Built-in fairness and explainability tools (Vertex AI Explainable AI) ensure models are transparent and equitable.	<b>Increased customer trust</b> and <b>reduced reputational risk</b> associated with biased or opaque models.
<b>Data Security</b>	Keeps all sensitive data and model artifacts within the secure GCP perimeter using VPC Service Controls.	<b>Elimination of third-party data exposure risk.</b> <b>Reduced cost</b> of data breach response and remediation.
<b>Cost Optimization</b>	Efficient resource utilization through managed services and autoscaling.	<b>15-25% lower compute costs</b> compared to self-managed MLOps infrastructure.

## Risk Mitigation Strategy

The architecture is specifically designed to mitigate key enterprise risks through the following mechanisms:

- **Data Leakage Prevention:** Strict **VPC Service Controls** and fine-grained **Cloud IAM** policies enforce a security perimeter, preventing unauthorized data movement.
- **Model Integrity:** Integrated **Vertex AI Model Monitoring** detects and alerts on model drift, data drift, and feature skew, ensuring model performance and reliability in production.
- **Ethical Compliance:** **Vertex AI Monitoring** and **Explainable AI** are used to ensure model fairness and detect bias, providing the necessary documentation for ethical review.
- **Auditability:** Comprehensive **Cloud Logging** and **Audit Logs** maintain immutable records of all API calls, resource changes, and pipeline executions, providing a complete forensic trail.

### 3. GRC Mapping (Governance, Risk, and Compliance)

---

This project adopts a **Security-by-Design** and **Compliance-by-Default** philosophy, aligning the MLOps implementation with leading global AI governance and security frameworks.

#### Compliance Frameworks Alignment

The project's design directly addresses the requirements of major AI and security frameworks:

Framework	Alignment Focus	Project Implementation
<b>NIST AI Risk Management Framework (AI RMF)</b>	Comprehensive AI risk management lifecycle, from design to deployment.	<b>Govern:</b> Use Model Cards for documentation. <b>Map:</b> Identify risks in the MLOps pipeline. <b>Measure:</b> Use Vertex AI Monitoring for performance and bias. <b>Manage:</b> Implement security controls (VPC-SC, IAM).
<b>ISO/IEC 42001:2023 (AI Management System)</b>	Standard for establishing, implementing, maintaining, and continually improving an AI Management System (AIMS).	<b>Systematic Governance:</b> Defined roles, responsibilities, and documented MLOps processes (Vertex AI Pipelines). <b>Risk Assessment:</b> Integrated risk analysis for data and model security.
<b>AI Act (EU)</b>	Requirements for high-risk AI systems (transparency, robustness, human oversight).	<b>Transparency:</b> Model Cards and Explainable AI. <b>Robustness:</b> Model Monitoring for drift/skew. <b>Oversight:</b> Defined human review points in the pipeline.
<b>GDPR (Article 22 &amp; 35)</b>	Automated decision-making transparency and Data Protection Impact Assessment (DPIA).	<b>Transparency:</b> Explainable AI provides insight into automated decisions. <b>DPIA Support:</b> Comprehensive logging and data access controls support DPIA requirements.
<b>SOC 2 (Trust Services Criteria)</b>	Security, Availability, Processing Integrity, Confidentiality, and Privacy.	<b>Security (CC6.1):</b> Strict IAM and VPC Service Controls. <b>Monitoring (CC7.2):</b> Continuous monitoring of production models and data access via Vertex AI Monitoring and Cloud Logging.

## Security Controls Implemented via Vertex AI

The following table details the specific GCP services and features used to enforce security and compliance controls:

Control	GCP Service	Purpose in PRJ-GCP-AI-082
<b>Explainable AI</b>	Vertex AI Explainable AI	Provides feature attributions for model predictions, ensuring transparency and supporting regulatory requirements for justification.
<b>Model Monitoring</b>	Vertex AI Model Monitoring	Continuously detects data drift, model drift, and feature skew in production, ensuring model reliability and compliance with robustness requirements.
<b>Data Encryption</b>	Cloud Storage, BigQuery	Data at rest and in transit is encrypted using Customer-Managed Encryption Keys (CMEK) or Google-Managed Encryption Keys (GMEK).
<b>Access Controls</b>	Cloud IAM	Enforces the <b>Principle of Least Privilege</b> across all resources, using fine-grained roles and conditional access policies.
<b>Bias Detection</b>	Vertex AI Monitoring	Proactively identifies and alerts on potential model bias across different feature slices, supporting fairness requirements.
<b>Security Perimeter</b>	VPC Service Controls (VPC-SC)	Creates a security boundary around Vertex AI, Cloud Storage, and BigQuery to prevent data exfiltration.
<b>Audit Trails</b>	Cloud Logging & Audit Logs	Captures all API calls, resource changes, and pipeline execution details for forensic analysis and immutable audit evidence.

## Audit Evidence Generation

The MLOps framework is configured to automatically generate and maintain critical artifacts required for audit and governance purposes:

- **Model Cards and Documentation:** Detailed records of model lineage, training data characteristics, intended use, and performance metrics are stored in the Vertex AI Model Registry.
- **Explainability Reports:** Outputs from Vertex AI Explainable AI for specific predictions are logged, providing justification for model decisions.

- **Bias Evaluation Results:** Metrics and reports from continuous model monitoring and fairness assessments are retained.
- **ML Pipeline Execution Logs:** Immutable records of every pipeline run, including input parameters, resource usage, and component outputs, are stored in Cloud Logging and GCS.

## 4. Prerequisites

---

Successful deployment requires a pre-configured GCP environment and local development tools.

### Required Accounts, Tools, and Permissions

1. **GCP Project:** A dedicated GCP project must be created, and billing must be enabled. The project ID is referred to as `$PROJECT_ID`.
2. **GCP CLI:** The `gcloud` and `gsutil` command-line tools must be installed, configured, and authenticated.
3. **Docker:** Docker must be installed locally for building the custom training container image.
4. **Permissions:** The user executing the deployment must have the following IAM roles on the project:
  - `roles/owner` or `roles/editor` (for initial setup).
  - `roles/serviceusage.serviceUsageAdmin` (to enable APIs).
  - `roles/iam.serviceAccountAdmin` (to create service accounts).
  - `roles/aiplatform.admin` (for full Vertex AI resource management).

### Initial Setup Commands

Execute the following commands to set up the environment variables, authenticate, and enable the necessary GCP APIs.

```
# 1. Define Environment Variables
export PROJECT_ID="PRJ-GCP-AI-082"
export REGION="us-central1" # Recommended region for Vertex AI services

# 2. Authenticate and set the active project
echo "Authenticating gcloud and setting project to $PROJECT_ID..."
gcloud auth login
gcloud config set project $PROJECT_ID

# 3. Enable Required GCP APIs
echo "Enabling required GCP APIs for Vertex AI, Storage, and IAM..."
gcloud services enable \
  cloudresourcemanager.googleapis.com \
  iam.googleapis.com \
  storage.googleapis.com \
  compute.googleapis.com \
  aiplatform.googleapis.com \
  artifactregistry.googleapis.com # Required for Docker image hosting
```

## 5. Architecture Overview

---

The solution is centered around a secure, containerized MLOps pipeline orchestrated by Vertex AI. The architecture adheres to the **Four Cornerstones of a Secure AI Platform**—Infrastructure, Data, Security, and Responsible AI—to ensure a holistic approach to governance and risk management [1].

## Key Architectural Components

Component	Role in MLOps Workflow	Security Consideration
<b>VPC Network</b>	Foundation for all GCP resources, providing network isolation.	All Vertex AI resources (Endpoints, Training) are configured to use private endpoints (PSC) within the VPC to prevent public internet exposure.
<b>Cloud Storage (GCS)</b>	Stores raw data, processed data, model artifacts, and pipeline root files.	Encrypted at rest (CMEK/GMEK) and protected by VPC Service Controls perimeter.
<b>Artifact Registry</b>	Hosts the custom Docker container image used for model training.	Access controlled by IAM and secured by the VPC-SC perimeter.
<b>Vertex AI Pipelines</b>	Orchestrates the end-to-end ML workflow (data processing, training, evaluation, deployment).	Runs using a dedicated Service Account with least-privilege access.
<b>Vertex AI Training</b>	Executes the custom training job using the containerized code.	Isolated execution environment, leveraging the secure container image from Artifact Registry.
<b>Vertex AI Model Registry</b>	Centralized repository for versioning, metadata, and Model Cards.	Enforces governance by requiring Model Cards before deployment.
<b>Vertex AI Endpoint</b>	Hosts the trained model for online, low-latency predictions.	Secured with IAM, rate limiting, and optionally Cloud Armor/Load Balancer for external access.
<b>Vertex AI Model Monitoring</b>	Provides continuous post-deployment monitoring.	Detects and alerts on data drift, model drift, and feature skew, crucial for maintaining compliance and performance.

## System Context and Flow

- 1. Data Ingestion:** Training data is securely uploaded to a GCS bucket, which is protected by VPC Service Controls.

2. **Containerization:** The custom training code is packaged into a Docker image and pushed to Artifact Registry.
3. **Orchestration:** A Data Scientist triggers a Vertex AI Pipeline run. The pipeline uses the secure container image and the GCS data.
4. **Training:** The pipeline executes a custom training job on Vertex AI Training, which is isolated within the VPC perimeter.
5. **Governance & Deployment:** Upon successful training, the model is registered in the Vertex AI Model Registry, complete with a Model Card. If evaluation metrics meet the threshold, the model is deployed to a Vertex AI Endpoint.
6. **Continuous Monitoring:** Vertex AI Model Monitoring is configured on the Endpoint to continuously check for drift and bias in production traffic.

*(Note: A detailed architecture diagram would typically be inserted here, illustrating the flow and security boundaries.)*

## 6. Step-by-Step Implementation

---

This section provides the detailed, actionable steps to deploy the secure MLOps framework. It assumes the existence of a custom training script ( `trainer/task.py` ), a `Dockerfile` , and a compiled pipeline definition ( `pipeline.json` ).

### Step 6.1: Setup Cloud Storage and Artifact Registry

We create the necessary storage for data and model artifacts, and a private repository for the Docker image.

```

# Define bucket and repository names
export BUCKET_NAME="{PROJECT_ID}-ml-artifacts"
export REPO_NAME="vertex-ai-repo"

# 1. Create GCS bucket for data and artifacts
echo "Creating GCS bucket: gs://$BUCKET_NAME/"
gsutil mb -l $REGION gs://$BUCKET_NAME/

# 2. Create Artifact Registry repository for Docker images
echo "Creating Docker repository in Artifact Registry: $REPO_NAME"
gcloud artifacts repositories create $REPO_NAME \
  --repository-format=docker \
  --location=$REGION \
  --description="Docker repository for Vertex AI custom training images"

# 3. Configure Docker to use gcloud as a credential helper
echo "Configuring Docker for gcloud authentication..."
gcloud auth configure-docker $REGION-docker.pkg.dev

```

## Step 6.2: Prepare and Upload Training Data

Upload the training data to the newly created GCS bucket.

```

# Assuming your training data is in a local directory named 'data'
# Create a placeholder file for demonstration if it doesn't exist
mkdir -p data
echo "feature_1,feature_2,target" > data/training_data.csv
echo "10.5,2.1,0" >> data/training_data.csv
echo "11.2,3.5,1" >> data/training_data.csv

# Upload the training data to GCS
echo "Uploading training data to GCS..."
gsutil cp ./data/training_data.csv gs://$BUCKET_NAME/data/

```

## Step 6.3: Build and Push Custom Training Container

The custom training code is containerized for execution on Vertex AI Training.

```
# Define the full image URI
export IMAGE_URI="$REGION-
docker.pkg.dev/$PROJECT_ID/$REPO_NAME/custom_trainer:latest"

# 1. Build the Docker image (assuming Dockerfile is in the current
directory)
# The Dockerfile should include your training script and dependencies.
echo "Building Docker image: $IMAGE_URI"
docker build -t $IMAGE_URI .

# 2. Push the image to Artifact Registry
echo "Pushing image to Artifact Registry..."
docker push $IMAGE_URI
```

## Step 6.4: Configure and Run Vertex AI Pipeline

The pipeline orchestrates the entire MLOps workflow. This step requires setting up a dedicated Service Account and defining pipeline parameters.

### 6.4.1: Create and Configure Pipeline Service Account

The Service Account (SA) needs permissions to read data from GCS, run Vertex AI jobs, and write artifacts.

```

export PIPELINE_SA_NAME="vertex-pipeline-sa"
export
PIPELINE_SA="${PIPELINE_SA_NAME}@${PROJECT_ID}.iam.gserviceaccount.com"

# 1. Create the service account (if it doesn't exist)
echo "Creating Vertex AI Pipeline Service Account: $PIPELINE_SA"
gcloud iam service-accounts create $PIPELINE_SA_NAME \
  --display-name "Vertex AI Pipeline Service Account for PRJ-GCP-AI-082"

# 2. Grant necessary roles to the SA
# roles/storage.objectAdmin: To read/write data and artifacts in GCS
# roles/aiplatform.user: To run training jobs and manage models/endpoints
echo "Granting IAM roles to the Service Account..."
gcloud projects add-iam-policy-binding $PROJECT_ID \
  --member="serviceAccount:$PIPELINE_SA" \
  --role="roles/storage.objectAdmin"

gcloud projects add-iam-policy-binding $PROJECT_ID \
  --member="serviceAccount:$PIPELINE_SA" \
  --role="roles/aiplatform.user"

```

## 6.4.2: Define Pipeline Parameters

The pipeline run requires a parameter file ( `pipeline_params.json` ) to define runtime configurations.

```

// pipeline_params.json
{
  "project_id": "PRJ-GCP-AI-082",
  "region": "us-central1",
  "gcs_data_path": "gs://PRJ-GCP-AI-082-ml-
artifacts/data/training_data.csv",
  "model_display_name": "Custom_Model_PRJ_082",
  "container_uri": "us-central1-docker.pkg.dev/PRJ-GCP-AI-082/vertex-ai-
repo/custom_trainer:latest",
  "machine_type": "n1-standard-4",
  "model_threshold": 0.75 // Example: Minimum acceptable evaluation metric
}

```

### 6.4.3: Execute the Pipeline Run

Assuming the pipeline definition (`pipeline.json`) has been compiled from a Kubeflow Pipelines SDK definition, execute the run.

```
# Run the pipeline
echo "Starting Vertex AI Pipeline run..."
gcloud ai pipelines runs create \
  --region=$REGION \
  --display-name="custom-training-run-$(date +%Y%m%d%H%M%S)" \
  --pipeline-root="gs://$BUCKET_NAME/pipeline_root" \
  --pipeline-file="./pipeline.json" \
  --parameter-file="./pipeline_params.json" \
  --service-account=$PIPELINE_SA
```

## 7. Validation & Testing

---

After deployment, rigorous validation is essential to confirm the MLOps pipeline is functioning correctly and the deployed model is accessible.

### 7.1: Pipeline Execution Validation

Monitor the pipeline run to ensure all components—data processing, custom training, model evaluation, and deployment—complete successfully.

- 1. Monitor Status in Console:** The most detailed view is in the GCP Console under **Vertex AI -> Pipelines**.
- 2. Check Logs via CLI:** Use the `gcloud` CLI to track the status of the latest run.

```
# Get the latest pipeline run ID
export LATEST_RUN_ID=$(gcloud ai pipelines runs list --region=$REGION --
format="value(name)" --limit=1 --sort-by=~createTime")

# Stream the state of the run (will show STATE_RUNNING, STATE_SUCCEEDED,
etc.)
echo "Monitoring pipeline run: $LATEST_RUN_ID"
gcloud ai pipelines runs describe $LATEST_RUN_ID --region=$REGION --
format="value(state)"

# Once the state is SUCCEEDED, verify the model is in the Model Registry.
```

## 7.2: Model Endpoint Prediction Test

Once the model is deployed to the Vertex AI Endpoint, a sample prediction request must be sent to validate its functionality and latency.

1. **Identify the Endpoint:** Retrieve the ID of the deployed endpoint.
2. **Prepare Sample Input:** Create a JSON file ( `sample_input.json` ) with data formatted exactly as the model expects.

```
// sample_input.json (Example for a model expecting two features)
{
  "instances": [
    {"feature_1": 15.0, "feature_2": 4.0},
    {"feature_1": 9.8, "feature_2": 1.5}
  ]
}
```

1. **Send Prediction Request:**

```
# Get the deployed model endpoint ID
export ENDPOINT_ID=$(gcloud ai endpoints list --region=$REGION --
filter="display_name:Custom_Model_PRJ_082" --format="value(name)")

# Send the prediction request
echo "Sending prediction request to Endpoint ID: $ENDPOINT_ID"
gcloud ai endpoints predict $ENDPOINT_ID \
  --region=$REGION \
  --json-request="sample_input.json" \
  --format="json"
```

**Expected Output:** The command should return a JSON response containing the model's predictions for the provided instances, confirming the model is live and serving traffic.

## 8. Troubleshooting

---

Common issues encountered during the deployment and operation of this MLOps framework, along with their resolutions.

Issue	Potential Cause	Resolution
<b>Pipeline Fails with Permission Denied</b>	The Pipeline Service Account ( <code>vertex-pipeline-sa</code> ) lacks necessary IAM roles for a specific resource (e.g., reading data, writing to Artifact Registry).	<b>Grant the SA the required roles.</b> Ensure it has <code>roles/storage.objectAdmin</code> for the GCS bucket and <code>roles/aiplatform.user</code> for Vertex AI operations. Check Cloud Audit Logs for the exact permission that was denied.
<b>Custom Training Job Fails</b>	Error in the training script ( <code>trainer/task.py</code> ), missing dependencies in the Docker image, or insufficient machine resources.	<b>Check the detailed logs in Cloud Logging</b> for the custom job. Verify the <code>Dockerfile</code> correctly installs all dependencies (e.g., in <code>requirements.txt</code> ). If the job runs out of memory, increase the <code>machine_type</code> in <code>pipeline_params.json</code> .
<b>Prediction Fails on Endpoint</b>	Input format in <code>sample_input.json</code> does not match the model's expected input signature, or the model is not fully deployed.	<b>Verify the model's input signature</b> in the Vertex AI Model Registry. Ensure the endpoint status is <code>ACTIVE</code> . If the model was recently deployed, wait a few minutes for the serving containers to initialize.
<b>Docker Push Fails (Permission Denied)</b>	Docker is not correctly authenticated to Artifact Registry, or the user lacks the <code>roles/artifactregistry.writer</code> role.	<b>Re-run <code>gcloud auth configure-docker \$REGION-docker.pkg.dev</code></b> . Ensure the user account has the necessary Artifact Registry permissions.
<b>VPC-SC Perimeter Violation</b>	A service account or resource attempted to access a restricted service from outside the perimeter.	<b>Review the VPC Service Controls logs</b> in Cloud Logging. Add the violating service account or IP range to the perimeter's access level, or ensure the resource is correctly placed within the protected project.

## 9. Cost Optimization

---

Managing costs is a critical aspect of MLOps in the cloud. The following strategies ensure efficient resource utilization on GCP.

### 1. Managed Compute Services:

- **Vertex AI Managed Training:** Always use Vertex AI Managed Training instead of self-managed Compute Engine instances. Vertex AI automatically provisions and **scales down compute resources immediately after the job completes**, eliminating idle compute costs.
- **Preemptible VMs:** For non-critical or batch training jobs, utilize preemptible virtual machines (VMs) to achieve significant cost savings (up to 80% off standard prices).

### 2. Machine Type Selection:

- **Start Small, Scale Up:** Begin with smaller, cost-effective machine types (e.g., `n1-standard-4` or `e2-standard-4`). Only scale up to larger machines or GPUs if performance benchmarks demonstrate a clear need and a positive ROI on the increased cost.
- **Custom Machine Types:** Use custom machine types to precisely match CPU and memory to the workload, avoiding the cost of over-provisioned standard types.

### 3. Storage Lifecycle Management:

- **GCS Object Lifecycle:** Implement Object Lifecycle Management policies on the GCS artifact bucket ( `$BUCKET_NAME` ). Automatically transition older model versions, raw logs, and intermediate pipeline artifacts to cheaper storage classes (Nearline or Coldline) after a defined period (e.g., 30 or 90 days).
- **BigQuery Partitioning:** Ensure BigQuery datasets are partitioned and clustered to minimize the amount of data scanned during queries, reducing analysis costs.

### 4. Endpoint Scaling and Idle Management:

- **Minimum Replica Count:** Configure the Vertex AI Endpoint with the **minimum replica count** set to 0 or 1. Setting it to 0 allows the endpoint to scale down completely during periods of zero traffic, minimizing idle costs.
- **Autoscaling:** Configure aggressive autoscaling policies to handle peak load efficiently while ensuring the minimum replica count keeps costs low during off-peak hours.

## 10. Security Best Practices

---

Security is the foundational pillar of this project. Beyond the initial setup, continuous adherence to best practices is mandatory for maintaining a hardened MLOps environment.

### 10.1: Network Perimeter Defense (VPC Service Controls)

For high-risk or highly regulated environments, **VPC Service Controls (VPC-SC)** are mandatory. VPC-SC creates a security perimeter around Google Cloud services, preventing data exfiltration by blocking unauthorized access from outside the perimeter.

**Implementation Note:** VPC-SC configuration is complex and typically requires the `Organization Admin` role. The following is a high-level example:

```
# Example: Create a service perimeter (requires Organization Admin role)
# This perimeter restricts access to Vertex AI, Storage, and BigQuery
# to resources ONLY within the project PRJ-GCP-AI-082.
echo "Creating VPC Service Controls perimeter (requires elevated
permissions)..."
gcloud access-context-manager perimeters create "prj_082_perimeter" \
  --perimeter-type="regular" \
  --resources="projects/$PROJECT_ID" \
  --restricted-
services="aiplatform.googleapis.com,storage.googleapis.com,bigquery.googleapis
\
  --description="Perimeter for PRJ-GCP-AI-082 MLOps resources"
```

## 10.2: Principle of Least Privilege (IAM)

All Service Accounts (SAs) and user accounts must be granted only the minimum necessary IAM roles to perform their specific function. Custom roles should be used where standard roles grant excessive permissions.

Persona / Service Account	Required Roles (Example)	Rationale
Data Scientist (User)	<code>roles/aiplatform.user</code> , <code>roles/storage.objectViewer</code>	Can run jobs and view data, but cannot modify infrastructure or delete artifacts.
Pipeline SA	<code>roles/aiplatform.serviceAgent</code> , <code>roles/storage.objectAdmin</code> (for artifact bucket only)	Can execute the pipeline and manage artifacts, but cannot manage user accounts or project billing.
Model Deployer (User)	<code>roles/aiplatform.endpointAdmin</code>	Can manage endpoints, but cannot train models or access raw data.

## 10.3: Secure MLOps Workflow Stages [1]

Security must be applied across the entire MLOps lifecycle:

### 1. Development and Data Ingestion:

- **Isolation:** Use private endpoints (PSC) for Vertex AI Notebooks to prevent direct internet access.
- **Sanitization:** Implement data validation and sanitization steps in the pipeline to prevent injection attacks and ensure data integrity.

### 2. Code and Pipeline Security:

- **Secure Repositories:** Use Cloud Source Repositories with IAM and branch protection policies.
- **Vulnerability Scanning:** Conduct vulnerability scanning on all custom training container images before pushing to Artifact Registry.

### 3. Training and Model Protection:

- **Confidential Computing:** For highly sensitive data, use **Confidential Computing** on Vertex AI to encrypt VM memory and data in use, even from Google operators.
- **Private Endpoints:** Ensure training jobs use private endpoints to access data and services.

### 4. Deployment and Serving:

- **Authentication:** Secure model endpoints with strong IAM policies.
- **Rate Limiting:** Implement rate limiting via Cloud Armor/Load Balancer to prevent abuse and denial-of-service attacks on the prediction endpoint.

### 5. Monitoring and Governance:

- **Continuous Monitoring:** Use Vertex Model Monitoring to continuously check for security-relevant anomalies, such as sudden changes in prediction distribution that could indicate a data poisoning attack.
- **Data Encryption:** Verify that the GCS bucket and BigQuery datasets are configured with **Customer-Managed Encryption Keys (CMEK)** for maximum control over encryption keys.

## 11. Cleanup

---

To avoid incurring future costs, it is essential to remove all resources created for this project.

```
# 1. Define Environment Variables (if not already set)
export PROJECT_ID="PRJ-GCP-AI-082"
export REGION="us-central1"
export BUCKET_NAME="${PROJECT_ID}-ml-artifacts"
export REPO_NAME="vertex-ai-repo"
export PIPELINE_SA="vertex-pipeline-sa@$PROJECT_ID.iam.gserviceaccount.com"

# 2. Undeploy and delete the Vertex AI Endpoint
echo "Undeploying and deleting the Vertex AI Endpoint..."
export ENDPOINT_ID=$(gcloud ai endpoints list --region=$REGION --
filter="display_name:Custom_Model_PRJ_082" --format="value(name)")
if [ -n "$ENDPOINT_ID" ]; then
    gcloud ai endpoints delete $ENDPOINT_ID --region=$REGION --quiet
else
    echo "No endpoint found to delete."
fi

# 3. Delete the Artifact Registry repository (deletes all stored images)
echo "Deleting Artifact Registry repository: $REPO_NAME"
gcloud artifacts repositories delete $REPO_NAME --location=$REGION --quiet

# 4. Delete the GCS bucket (WARNING: This deletes all data and artifacts)
echo "Deleting GCS bucket: gs://$BUCKET_NAME"
gsutil rm -r gs://$BUCKET_NAME

# 5. Delete the Service Account
echo "Deleting Service Account: $PIPELINE_SA"
gcloud iam service-accounts delete $PIPELINE_SA --quiet

# 6. (Optional) Delete the entire GCP project
# WARNING: This is irreversible and deletes ALL resources in the project.
# gcloud projects delete $PROJECT_ID
```

---

## References

---

[1] Schneider Larbi and David Peterside, "Mastering secure AI on Google Cloud: A practical guide for enterprises," *Google Cloud Blog*, March 21, 2025. <https://cloud.google.com/blog/products/identity-security/mastering-secure-ai-on-google-cloud-a-practical-guide-for-enterprises>