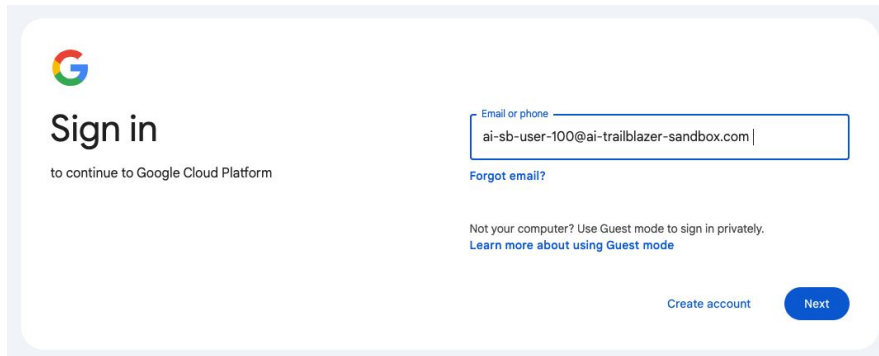# LAB GUIDE

AI Labs - Prudential with Google

# **Labs Access:** Accessing Google Cloud Console
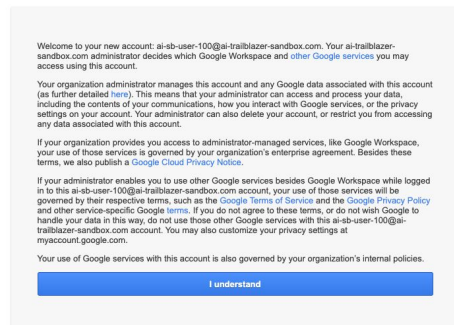
**Step 1:** Navigate to Google Cloud Console.
https://console.cloud.google.com/

**Step 2:** For the first login, you have to click "I Understand"

# Labs Access: Accessing Google Cloud Console

**Step 3:** You will be prompted to change password. Please enter a new password.

**Step 4:** Check the checkbox for "Terms of Service" and then proceed to click on "Agree and Continue".

# Labs Access: Accessing Google Cloud Console

Navigate to Vertex AI and then Workbench

# Labs Setup: Accessing Workbench Notebook

**Step 4:** Once at Workbench, ensure you are on the **INSTANCES** tab

**Step 5:** You should see that the notebooks created for your team.
Each user is assigned 2 notebooks. 1 with GPU and another
without GPU.

GPU instance: ai-gpu-100-XXX. Non-GPU instance: ai-no-gpu-100--XXX

# Labs Setup: Accessing Workbench Notebook

**Step 6:** Select the 2 instances assigned and click on **START** in the top menu. If the instances are already started, skip to step 8.

# Labs Setup: Accessing Workbench Notebook

**Step 7:** Wait 1-3 minutes for the instances to start. Verify that there is a green tick beside the instances.



**Step 8:** Click on **OPEN JUPYTERLAB** for the **gpu instance** <ai-gpu-100-XXX>.
A new tab will open with access to the Jupyter Notebook.



Google

# Labs Setup: Accessing Workbench Notebook

**Step 9**: Expand the **Git** menu and click on **Clone a Repository**

**Step 10**: Paste the URL provided for you into the text box and click **Clone**

URL = https://github.com/analyticsrepo01/pru-ai-labs.git



Google

# Labs Setup: Setting Up Notebooks

**Step 12**: Expand the **Run** menu and select **Run All Cells**

# Image model
with custom
training

# Labs Setup: Setting Up Notebooks

**Step 11**: Double click **pru_sd_xl_finetuning_dreambooth_lora.ipynb** to open the next notebook

# Labs Setup: Setting Up Notebooks

**Step 12**: Expand the **Run** menu and select **Run All Cells**.

# Open LLM
# with custom
# training

# Labs Setup: Tuning a Open source LLM

**Step 13**: On Google Cloud Console, Navigate to Vertex AI → Model Garden

TOOLS

Dashboard

Model Garden

Pipelines

NOTEBOOKS ⌄

Colab Enterprise

Workbench

VERTEX AI STUDIO ⌄

Overview

Freeform

Chat

Vision

Translation

Speech

Prompt gallery

Prompt management

Tuning

BUILD WITH GEN AI ⌄

Extensions

Marketplace

Model Garden   ✦ EXPLORE GENERATIVE AI   💡 VIEW MY ENDPOINTS & MODELS   ⚡ DEPLOY FROM HUGGING FACE   📑 VIEW RELEASE NOTES

Modalities

| | |
|---|---|
| Language | 67 |
| Vision | 88 |
| Tabular | 7 |
| Document | 8 |
| Speech | 2 |
| Video | 6 |
| Multimodal | 21 |
| Audio | 1 |

Tasks

| | |
|---|---|
| Generation | 74 |
| Classification | 66 |
| Detection | 44 |
| Extraction | 28 |
| Recognition | 26 |
| Translation | 24 |
| Embedding | 7 |
| Segmentation | 12 |
| Retrieval | 2 |
| Open vocabulary detection | 2 |
| Open vocabulary segmentation | 2 |

Q Search models    **Llama 3.1**

Browse, customize, and deploy machine learning models with **Model Garden**. Choose from models created by Google and other providers.

Gemini

Imagen 3
Generate images with text prompts

Gemma 2

Llama 3.2
Experience AI

Sort by:  **Trending**  Newest  Last Update

## Foundation models    → SHOW ALL (93)

Pre-trained multi-task models that can be further tuned or customized for specific tasks.

**Gemini 1.5 Pro**

Created from the ground up to be multimodal (text, images, videos) and to scale across a wide range of tasks

**Gemini 1.5 Flash**

The best performing Gemini model with features for a wide range of tasks

**Gemini 1.0 Pro**

Designed to balance quality, performance, and cost for tasks such as content generation, editing, summarization, and classification

**Gemini 1.0 Pro Vision**

Created to be multimodal (text, images, code) and to scale across a wide range of tasks

## Featured partners

ANTHROP\C

∞ Meta

🤗 Hugging Face

MISTRAL AI_

**Step 14**: Choose Llama 3.1 model, Click on "FINE-TUNE"

## Create a fine-tuned model

Fine tune this model using supervised tuning.

**Tuned model name ***
llama-3-1-70b-1731549015067

**Base model**
Llama-3-1-70B ▼ ❓

📁 Output directory *     BROWSE ❓
⚠ Input is required

**Region ***
us-central1 (Iowa) ▼

The base model will be tuned with the following settings ( See pricing ) ⬈

- Machine type: a2-ultragpu-8g
- Accelerator type: NVIDIA_A100_80GB
- Accelerator count: 8

## Tuning parameters

**Number of epochs ***
3     ❓

**Learning rate ***
0.0002     ❓

⚙ VIEW TUNING CONFIG

## Dataset

**Tuning dataset ***
✅ cloud-samples-data/vertex-ai/model-evaluation/peft_train_sam   BROWSE ❓

**Evaluation dataset ***
✅ cloud-samples-data/vertex-ai/model-evaluation/peft_eval_sam   BROWSE ❓

ⓘ VIEW EXAMPLE FORMAT

**START TUNING**    CANCEL

---

Next start tining once all filled up

---

Create or select the bucket

## Create a bucket

• **Get Started**

Pick a **globally unique**, permanent name. Naming guidelines ⬈

llama3trainingv1|

Tip: Don't include any sensitive information

Optimize storage for data-intensive workloads ⌄

Labels (optional) ⌄

**CONTINUE**

• **Choose where to store your data**

**Location**: us (multiple regions in United States)
**Location type**: Multi-region

• **Choose a storage class for your data**

**Default storage class**: Standard

• **Choose how to control access to objects**

**Public access prevention**: On
**Access control**: Uniform

• **Choose how to protect object data**

**Soft delete policy**: Default
**Object versioning**: Disabled
**Bucket retention policy**: Disabled
**Object retention**: Disabled
**Encryption type**: Google-managed

**CREATE**    CANCEL

Google

Google Cloud | prusandbx-nprd-uat-u9pahg ▾

iam | ✕ | 🔍 Search

**Vertex AI** 📌

**Training** | ➕ **TRAIN NEW MODEL** | ↻

Go to training tab - from left menu & check training started?

**TRAINING PIPELINES** | **CUSTOM JOBS** | H

BUILD WITH GEN AI ⌄

- 🔌 Extensions
- 🔲 Code samples
- ➕ Agents

DATA ⌄

- 📦 Feature Store
- 🔲 Datasets

MODEL DEVELOPMENT ⌄

- ⦿ **Training**
- ⌛ Experiments
- ▦ Metadata

DEPLOY AND USE ⌄

- 📍 Model Registry
- ⦿ Online prediction
- 🗄 Batch predictions
- 📈 Monitoring
- ⁂ Vector Search

MANAGE ⌄

- ⇲ Ray on Vertex AI
- ⇥ Migrate to Vertex AI
- 🛒 Marketplace

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | pipeline | Custom | sec | 11:34:27 AM | 10:41:43 AM | AM |
| llama3-1-lora-train-20241201-164505 | 233993907094945792 | ✅ Finished | Training pipeline | ⊕ Custom | 17 min 6 sec | Dec 2, 2024, 1:02:14 AM | Dec 2, 2024, 12:45:07 AM | Dec 2, 2024, 1:02:14 AM | — | ⋮ |
| llama31-lora-20241201-143742 | 3612889896273838080 | ✅ Finished | Training pipeline | ⊕ Custom | 28 min 8 sec | Dec 1, 2024, 11:05:52 PM | Dec 1, 2024, 10:37:44 PM | Dec 1, 2024, 11:05:52 PM | — | ⋮ |
| llama3-1-lora-train-20241201-141259 | 3170974182838108160 | ❗ Failed | Training pipeline | ⊕ Custom | 52 min 45 sec | Dec 1, 2024, 11:05:45 PM | Dec 1, 2024, 10:13:00 PM | Dec 1, 2024, 11:05:45 PM | — | ⋮ |
| llama31-lora-20241201-143729-1733063850 | 2477982790176473088 | ✅ Finished | Training pipeline | ⊕ Custom | 28 min 8 sec | Dec 1, 2024, 11:05:39 PM | Dec 1, 2024, 10:37:30 PM | Dec 1, 2024, 11:05:39 PM | — | ⋮ |
| llama3-1-lora-train-20241201-144150 | 3286941873242898432 | ✅ Finished | Training pipeline | ⊕ Custom | 23 min 7 sec | Dec 1, 2024, 11:04:59 PM | Dec 1, 2024, 10:41:52 PM | Dec 1, 2024, 11:04:59 PM | — | ⋮ |
| llama3-1-lora-train-20241201-134300 | 3813018604715114496 | ❗ Failed | Training pipeline | ⊕ Custom | 52 min 14 sec | Dec 1, 2024, 10:35:16 PM | Dec 1, 2024, 9:43:01 PM | Dec 1, 2024, 10:35:16 PM | — | ⋮ |
| llama31-lora-20241201-143018 | 478384555623972864 | ❗ Failed | Training pipeline | ⊕ Custom | 1 min 32 sec | Dec 1, 2024, 10:31:52 PM | Dec 1, 2024, 10:30:20 PM | Dec 1, 2024, 10:31:52 PM | — | ⋮ |
| llama3-1-lora-train-20241201-141611 | 4672924658566168576 | ✅ Finished | Training pipeline | ⊕ Custom | 15 min 34 sec | Dec 1, 2024, 10:31:47 PM | Dec 1, 2024, 10:16:12 PM | Dec 1, 2024, 10:31:47 PM | — | ⋮ |
| llama3-1-lora-train-20241201-141539 | 5396878298665975808 | ✅ Finished | Training pipeline | ⊕ Custom | 15 min 35 sec | Dec 1, 2024, 10:31:16 PM | Dec 1, 2024, 10:15:40 PM | Dec 1, 2024, 10:31:16 PM | — | ⋮ |
| llama3-1-lora-train-20241201-141534 | 4221438795922276352 | ✅ Finished | Training pipeline | ⊕ Custom | 15 min 35 sec | Dec 1, 2024, 10:31:11 PM | Dec 1, 2024, 10:15:35 PM | Dec 1, 2024, 10:31:11 PM | — | ⋮ |
| llama3-1-lora-train-20241201-142214 | 1870559790434877440 | ❗ Failed | Training pipeline | ⊕ Custom | 5 min 2 sec | Dec 1, 2024, 10:27:19 PM | Dec 1, 2024, 10:22:16 PM | Dec 1, 2024, 10:27:19 PM | — | ⋮ |
| llama3-1-lora-train-20241201-141756 | 9148376788265598976 | ❗ Failed | Training pipeline | ⊕ Custom | 5 min 33 sec | Dec 1, 2024, 10:23:32 PM | Dec 1, 2024, 10:17:58 PM | Dec 1, 2024, 10:23:32 PM | — | ⋮ |
| llama3-1-lora-train-20241201-131514 | 6548955378342690816 | ❗ Failed | Training pipeline | ⊕ Custom | 53 min 14 sec | Dec 1, 2024, 10:08:31 PM | Dec 1, 2024, 9:15:16 PM | Dec 1, 2024, 10:08:31 PM | — | ⋮ |
| llama3-1-lora-train-20241201-134445 | 6964412443967619072 | ✅ Finished | Training pipeline | ⊕ Custom | 16 min 35 sec | Dec 1, 2024, 10:01:23 PM | Dec 1, 2024, 9:44:47 PM | Dec 1, 2024, 10:01:23 PM | — | ⋮ |
| llama3-1-lora-train-20241201-134438 | 3384050740208074752 | ✅ Finished | Training pipeline | ⊕ Custom | 16 min 5 sec | Dec 1, 2024, 10:00:46 PM | Dec 1, 2024, 9:44:40 PM | Dec 1, 2024, 10:00:46 PM | — | ⋮ |
| llama3-1-lora-train-20241201-133545 | 3630622810806600408 | ✅ Finished | Training | | 13 min 4 sec | Dec 1, 2024, 9:38:52 PM | Dec 1, 2024, 9:25:47 PM | Dec 1, 2024, 9:38:52 |

# Example on custom data: Llama3.1 on Insurance data

Open and run all the notebook pru_llama3_1_finetuning.ipynb

s-central1.notebooks.googleusercontent.com/lab/tree/GenAI8/pru-ai-labs/pru_llama3_1_finetuning.ipynb

Run ALL

n1-standard-4 ▾

Multi-agent_002.ipynb ✕ | Story_Multi_agent_002.i✕ | RAG_Crew.ipynb ✕ | pru_llama3_1_finetuning. ✕ | Terminal 1 ✕ | Terminal 2 ✕ +

💾 + ✂ ⎘ ⎘ ▶ ■ C ⏩ Markdown ▾ 🕐 git ⎘ Execute                                    ✺ PyTorch 1-13 (Local) ○ 🗓

/ GenAI8 / pru-ai-labs /

| Name | ▲ Last Modified |
|---|---|
| sd_models | 4 days ago |
| web-app | 5 days ago |
| video | 5 days ago |
| testing | 5 days ago |
| images | 5 days ago |
| backup_folder | 5 days ago |
| pru_llama3_1_finetuning.ipynb | seconds ago |
| pru_sd_xl_finetuning_dreambooth_lor... | 39 minutes ago |
| Multi-agent_002.ipynb | an hour ago |
| story.pdf | an hour ago |
| mdpdf.log | an hour ago |
| story.md | an hour ago |
| part6_Gemini_on_video.ipynb | 12 hours ago |
| pru_DIY_RAG_pdf.ipynb | 12 hours ago |
| pru_stable_diffusion_2_1.ipynb | 3 days ago |
| pru_RAG_VertexSearch.ipynb | 4 days ago |
| part_mixtral.ipynb | 4 days ago |
| requirements.txt | 5 days ago |
| questions_test.json | 5 days ago |
| questions.json | 5 days ago |
| part_llama3.ipynb | 5 days ago |
| part_Imagen.ipynb | 5 days ago |
| intro_multimodal_use_cases_latest.ip... | 5 days ago |
| pru_qa.csv | 5 days ago |
| Imagen_on_questions.ipynb | 5 days ago |
| Agent_Example.ipynb | 5 days ago |
| $BUCKET_URI | 5 days ago |
| Agent_Email_conversation002.ipynb | 5 days ago |

```
[20]:   1 # Copyright 2024 Google LLC
        2 #
        3 #
```

# Vertex AI Model Garden - Llama 3.1 Finetuning Insurance Data

CO
Run in Colab Enterprise   View on GitHub

▶ ⎘ ↑ ↓ ⬆ ⎯ 🗑 ⊕

## Overview

This notebook demonstrates finetuning and deploying Llama 3.1 models with Vertex AI. After finetuning, we can deploy models on Vertex with GPU.

### Objective

- Finetune Llama 3.1 models with Vertex AI Custom Training Jobs.
- Deploy finetuned Llama 3.1 models on Vertex AI Prediction.
- Send prediction requests to your finetuned Llama 3.1 models.

### Costs

This tutorial uses billable components of Google Cloud:

- Vertex AI
- Cloud Storage

## Before you begin

### Install dependencies

```
[21]:   1 print("Installing google-cloud-aiplatform")
```

# RAG for Prudential Use case

# Labs Setup: Accessing Workbench Notebook

**Step 15:** Click on **OPEN JUPYTERLAB** for the **no-gpu** instance <ai-no-gpu-100-XXX>. A new tab will open with access to the Jupyter Notebook.



Google

# Labs Setup: Accessing Workbench Notebook

**Step 16**: Expand the **Git** menu and click on **Clone a Repository**

**Step 17**: Paste the URL provided for you into the text box and click **Clone**

    URL = https://github.com/analyticsrepo01/pru-ai-labs.git



Google

# Labs Setup: Vertex AI Search

**Step 18:** Double click **pru_RAG_VertexSearch.ipynb** to open the Jupyter Notebook

**Step 19:** In the top menu bar, expand the **Run** menu and select **Run All Cells**

# DIY RAG

# **Labs Setup:** Vector AI Search

**Step 20:** Double click **pru_DIY_RAG_pdf.ipynb** to open the Jupyter Notebook

**Step 21:** In the top menu bar, expand the **Run** menu and select **Run All Cells**

# Kernel Restarts after installation

Step1 : Click on this cell after restart

Step2 : Run Selected Cell and All Below

# Multi Agents

ai-gpu-100-f20436e7

File   Edit   View   Run   Kernel   Git   Tabs   Settings   Help

g2-standard-4 ▾

Filter files by name

/ pru-ai-labs /
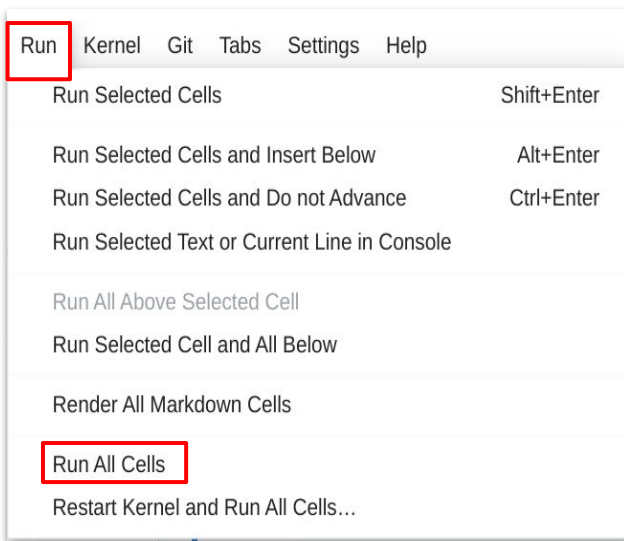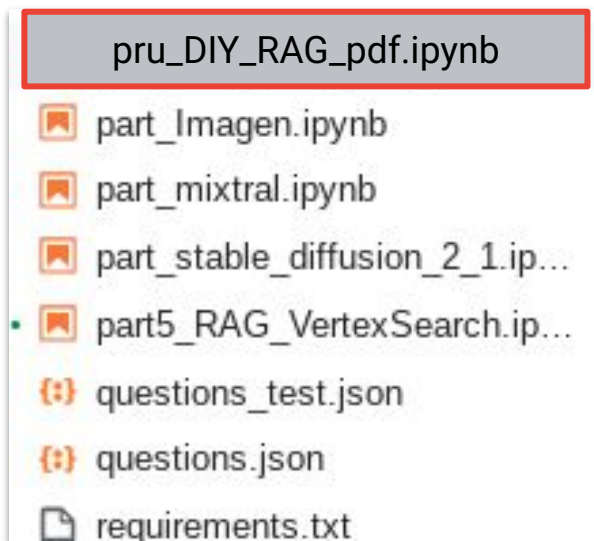
| Name | Last Modified |
|---|---|
| images | a day ago |
| results | a day ago |
| testing | a day ago |
| vertex-ai-samples | a day ago |
| video | a day ago |
| web-app | a day ago |
| $BUCKET_URI | a day ago |
| ● Agent_Email_conversation002.ip... | 2 minutes ago |
| Agent_Example.ipynb | a day ago |
| Imagen_on_questions.ipynb | a day ago |
| intro_multimodal_use_cases_lat... | a day ago |
| mdpdf.log | 6 minutes ago |
| Multi-agent_002.ipynb | seconds ago |
| part_DIY_RAG_pdf.ipynb | a day ago |
| part_Imagen.ipynb | a day ago |
| part_llama3.ipynb | a day ago |
| part_mixtral.ipynb | a day ago |
| part6_Gemini_on_video.ipynb | a day ago |
| ● pru_llama3_1_finetuning.ipynb | 4 hours ago |
| pru_qa.csv | a day ago |
| pru_RAG_VertexSearch.ipynb | a day ago |
| pru_sd_xl_finetuning_dreamboo... | a day ago |
| pru_stable_diffusion_2_1.ipynb | a day ago |
| {} questions_test.json | a day ago |
| {} questions.json | a day ago |

Terminal 1 ✕   pru_llama3_1_finetuning.ip ✕   Agent_Email_conversation ●   Multi-agent_002.ipynb ✕ +

Code ▾   git   Execute

PyTorch 1-13 (Local)

file_read_tool        1/2

```
[1]: # !conda create -n crewai python=3.11
     # !conda activate crewai -
     !pip install -q --upgrade google-cloud-aiplatform
     !pip install -q -U 'crewai[tools]' mdpdf
```

ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the so
urce of the following dependency conflicts.
gcsfs 2024.10.0 requires fsspec==2024.10.0, but you have fsspec 2024.9.0 which is incompatible.
kfp 2.5.0 requires kubernetes<27,>=8.0.0, but you have kubernetes 31.0.0 which is incompatible.
kfp 2.5.0 requires requests-toolbelt<1,>=0.8.0, but you have requests-toolbelt 1.0.0 which is incompatible.
kfp 2.5.0 requires urllib3<2.0.0, but you have urllib3 2.2.3 which is incompatible.

```
[2]: import re

     PROJECT_ID = !(gcloud config get-value core/project)
     PROJECT_ID = PROJECT_ID[0]

     SVC_ACC = !(gcloud config get-value core/account)
     SVC_ACC = SVC_ACC[0]

     PROJECT_NUMBER=str(re.search(r'\d+', SVC_ACC).group())

     LOCATION="asia-southeast1"

     FOLDER_NAME="."
```

```
[3]: from crewai import Agent, Task, Crew, Process
     from crewai_tools import tool
     from langchain_openai import ChatOpenAI
     from crewai_tools.tools import FileReadTool
     import os, requests, re, mdpdf, subprocess
     from openai import OpenAI
```

```
[4]: !pip install --upgrade --quiet  langchain-core langchain-google-vertexai
     !pip install mdpdf
```

ERROR: pip's dependency resolver does not currently take into account all the packages that are installed. This behaviour is the so
urce of the following dependency conflicts.

Google

# Optional Lab :

Go to this github URL and clone

https://github.com/analyticsrepo0
1/ai_agents_v2

git clone <URL>

# Labs Setup: Setting Up Notebooks

## Note:

**Avoid letting your laptops/computers enter sleep mode to prevent problems arising when running the notebooks**

# Thank you

# Multi-Agent Lab: Generating AI News Insights

**Objectives:**
- Retrieve news about AI using RAG
- Generate summarized reports on AI trends