# Production ML Pipelines
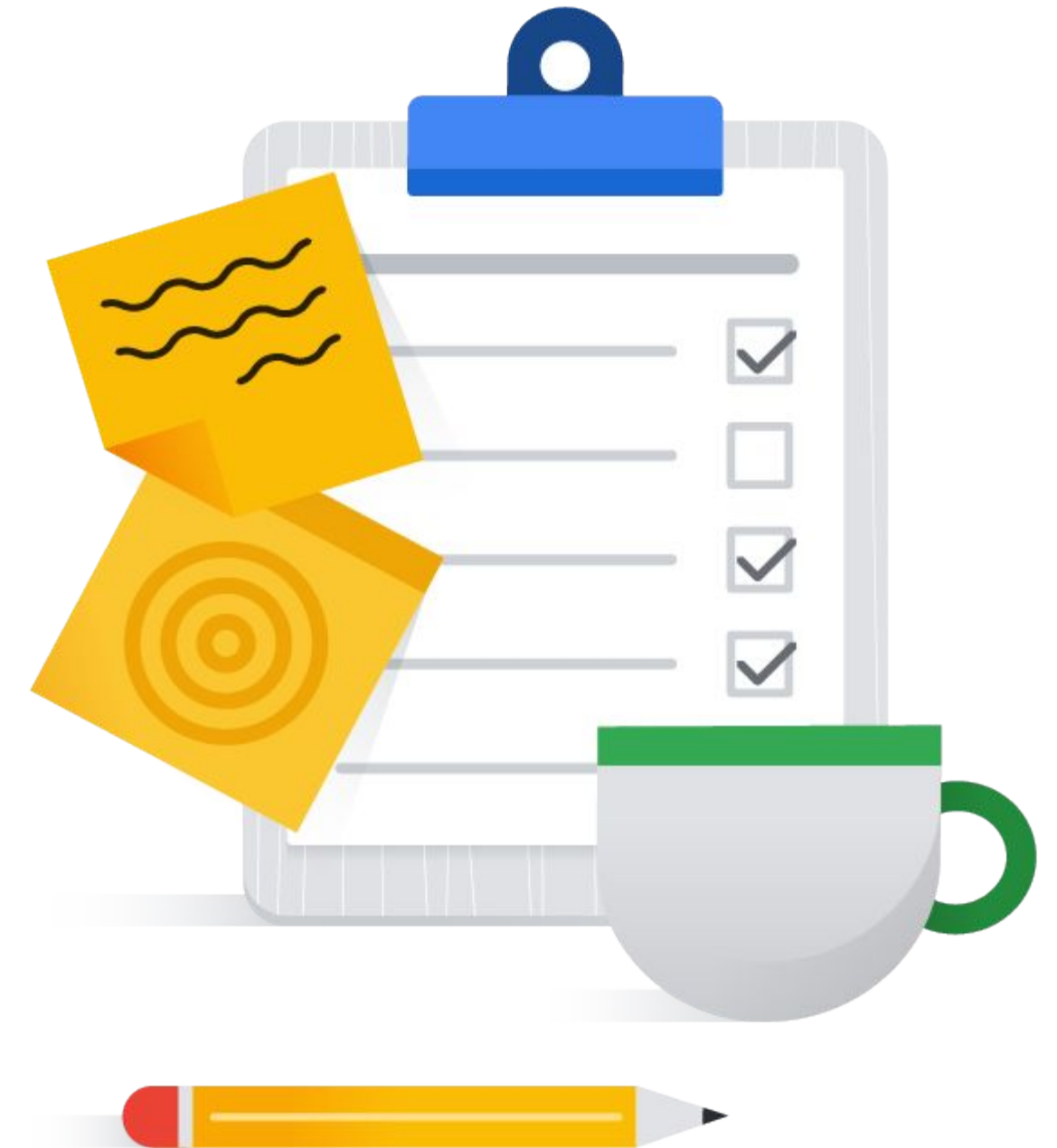
**04**

# Production ML Pipelines

| 01 | Ways to do ML on Google Cloud |
|----|-------------------------------|
| 02 | Vertex AI Pipelines |
| 03 | AI Hub |

Google Cloud

# Production ML Pipelines

| | |
|---|---|
| **01** | **Ways to do ML on Google Cloud** |
| **02** | Vertex AI Pipelines |
| **03** | AI Hub |

# Create and deploy custom models with Vertex AI

Cloud TPUs

Compute Engine

Dataproc

Google Kubernetes Engine

Vertex AI

BigQuery ML

**Build a Custom Model**

## AutoML

Cloud Translation API

Vision API

Speech-to-Text API

Video Intelligence API

Data Loss Prevention API

Text-to-Speech API

Cloud Natural Language API

Dialogflow

**Build Custom Model (codeless)**

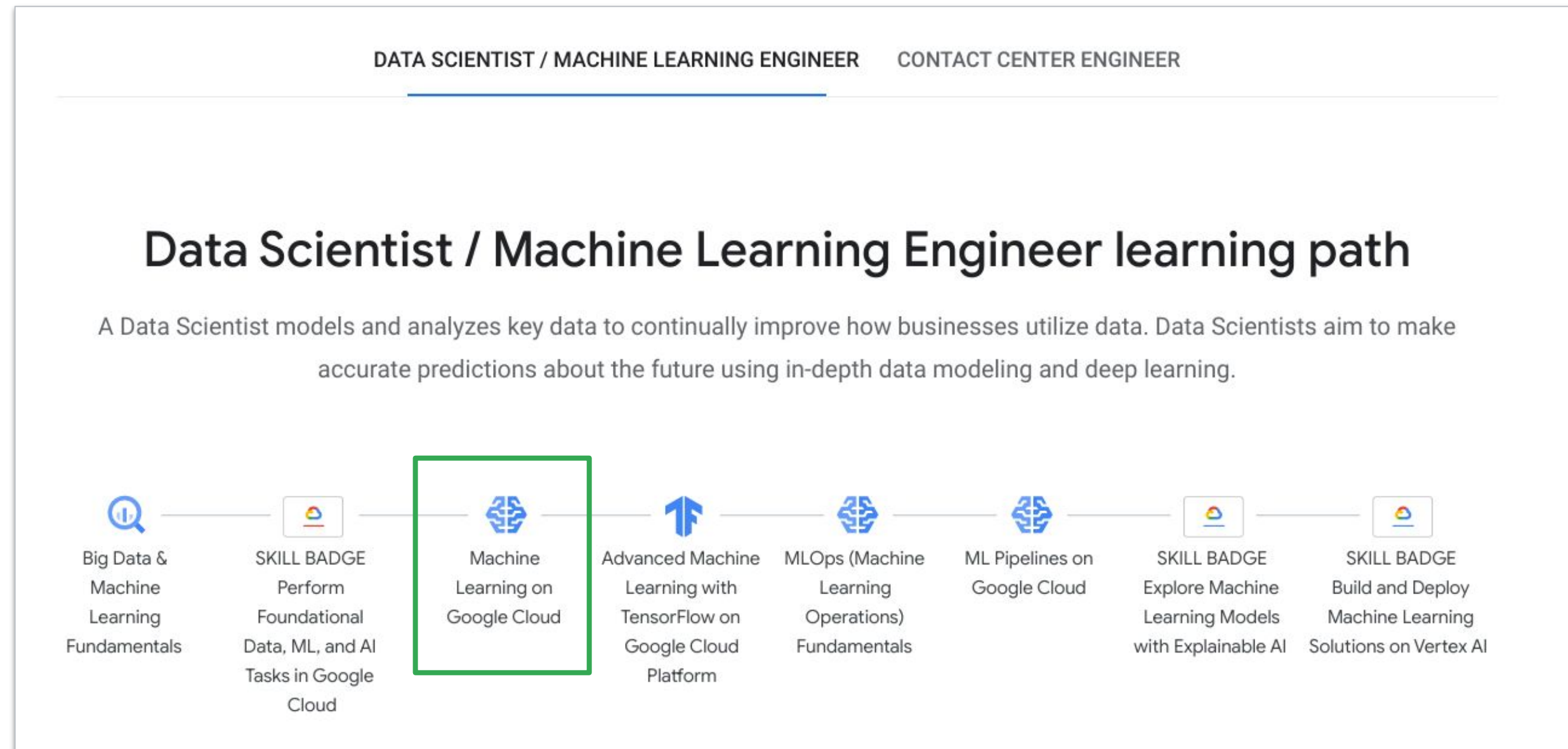**Call a Pretrained Model**

Google Cloud

# Vertex AI is a fully managed service for custom machine learning models



Vertex AI

Cloud TPU

- Scales to production
- Batching and distribution of model training
- Performs transformations on input data
- Hyper-parameter tuning
- Host and autoscale predictions
- Serverless - self-tuning - manages overhead
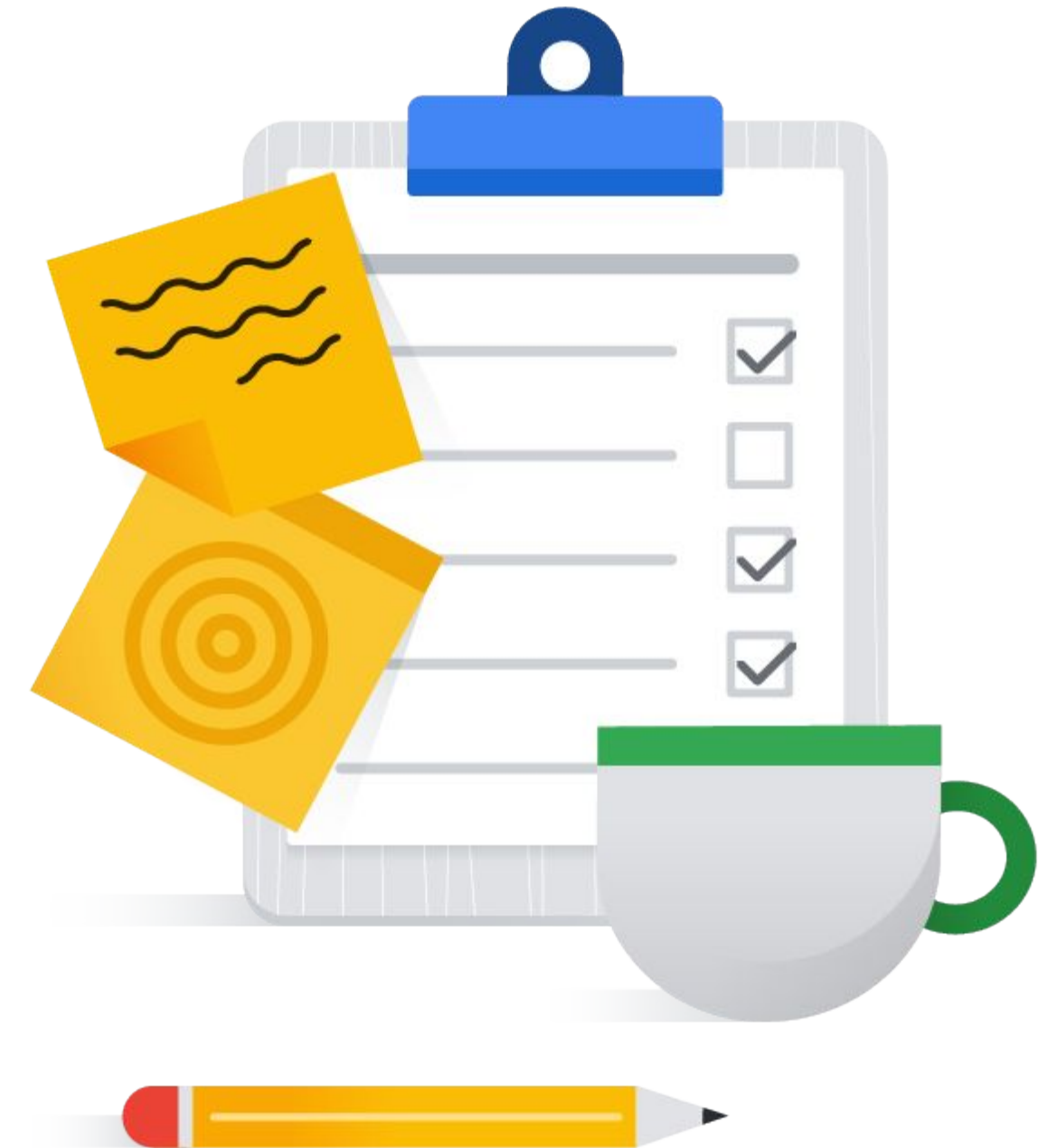
Google Cloud

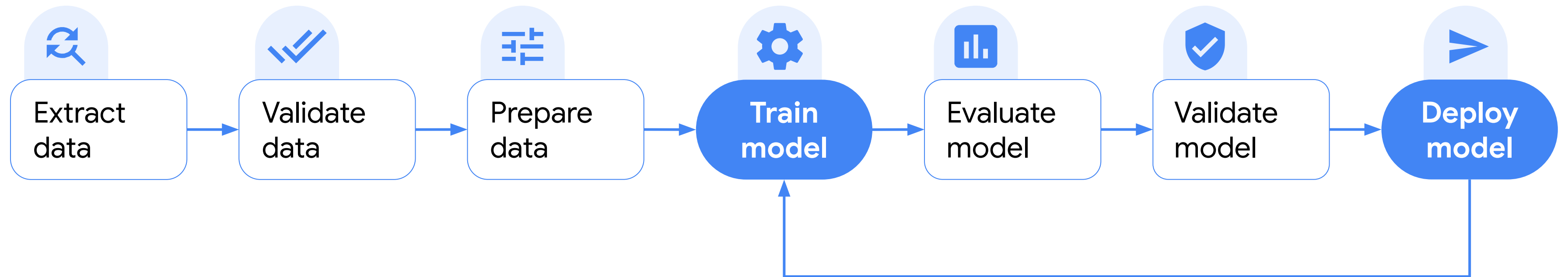# In this course, we don't cover writing TensorFlow models, only ways to operationalize them



Google Cloud Training - Machine Learning and AI

# Production ML Pipelines

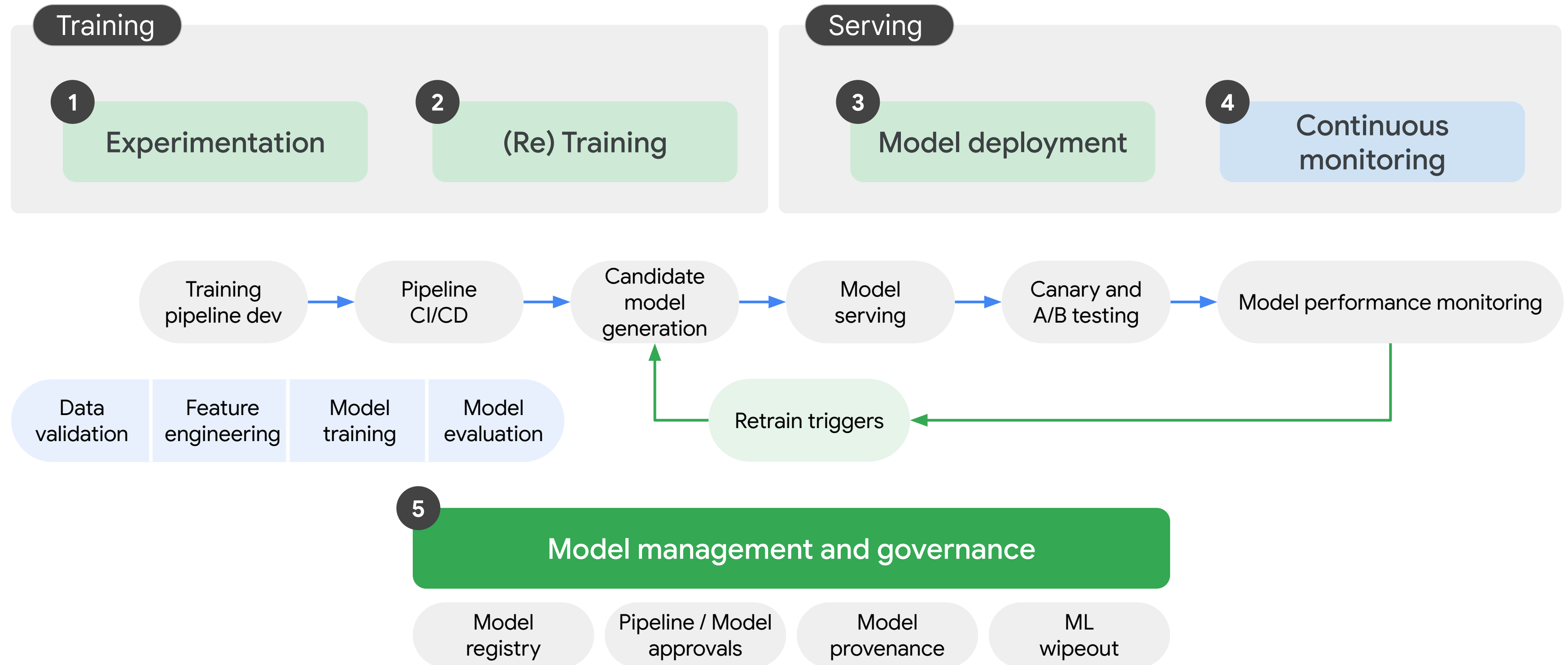| 01 | Ways to do ML on Google Cloud |
|----|-------------------------------|
| 02 | **Vertex AI Pipelines** |
| 03 | AI Hub |

# Pipelines automate the training and deployment of models



Extract data → Validate data → Prepare data → **Train model** → Evaluate model → Validate model → **Deploy model**

# Pipelines are the backbone of production ML systems

**Training**

**1** Experimentation

**2** (Re) Training

**Serving**

**3** Model deployment

**4** Continuous monitoring

Training pipeline dev → Pipeline CI/CD → Candidate model generation → Model serving → Canary and A/B testing → Model performance monitoring

Data validation | Feature engineering | Model training | Model evaluation

Retrain triggers

**5** Model management and governance

Model registry | Pipeline / Model approvals | Model provenance | ML wipeout

Google Cloud

# Pipelines product portfolio

### Kubeflow

**Kubeflow Pipelines**

- Kubernetes-native.
- Open source.
- The industry standard for running ML Pipelines.

### Google Cloud

**AI Platform Pipelines - Hosted** [Beta]

- Kubeflow pipelines running on Google Cloud.
- Optimized for GKE.
- Integrated with Google Cloud services.

**Vertex Pipelines - Managed** [PREVIEW]

- Fully managed and serverless.
- Allows users to focus on building their pipelines, scale easily, and pay only for the resources they use.

# Write your pipeline

## Easy to use Python SDKs

Build pipelines using Data Scientist friendly SDKs like TensorFlow Extended and Kubeflow Pipelines.

## Rich, scalable pre-built components

We provide a rich set of pre-built components for common ML tasks, which leverage Google Cloud services.

```python
@dsl.pipeline(pipeline_root=PIPELINE_ROOT, name="metadata-pipeline-v2")
def pipeline(message: str):
    importer = kfp.dsl.importer(
        artifact_uri="gs://ml-pipeline-playground/shakespeare1.txt",
        artifact_class=Dataset,
        reimport=False,
    )
    preprocess_task = preprocess(message=message)
    train_task = train(
        dataset_one=preprocess_task.outputs["output_dataset_one"],
        dataset_two=preprocess_task.outputs["output_dataset_two"],
        imported_dataset=importer.output,
        message=preprocess_task.outputs["output_parameter"],
        num_steps=5,
    )
    read_task = read_artifact_input(
        train_task.outputs["generic_artifact"]
    )
```

Google Cloud

# Key capabilities

## Python SDKs

Data scientist friendly Python SDKs

**1**

## Scalable

Run as many pipelines on as much data as you want.

**2**

## Metadata and lineage

Store metadata for every artifact produced by the pipeline.

**3**

## Monitoring UIs and APIs

Track and debug pipelines executions.

**4**

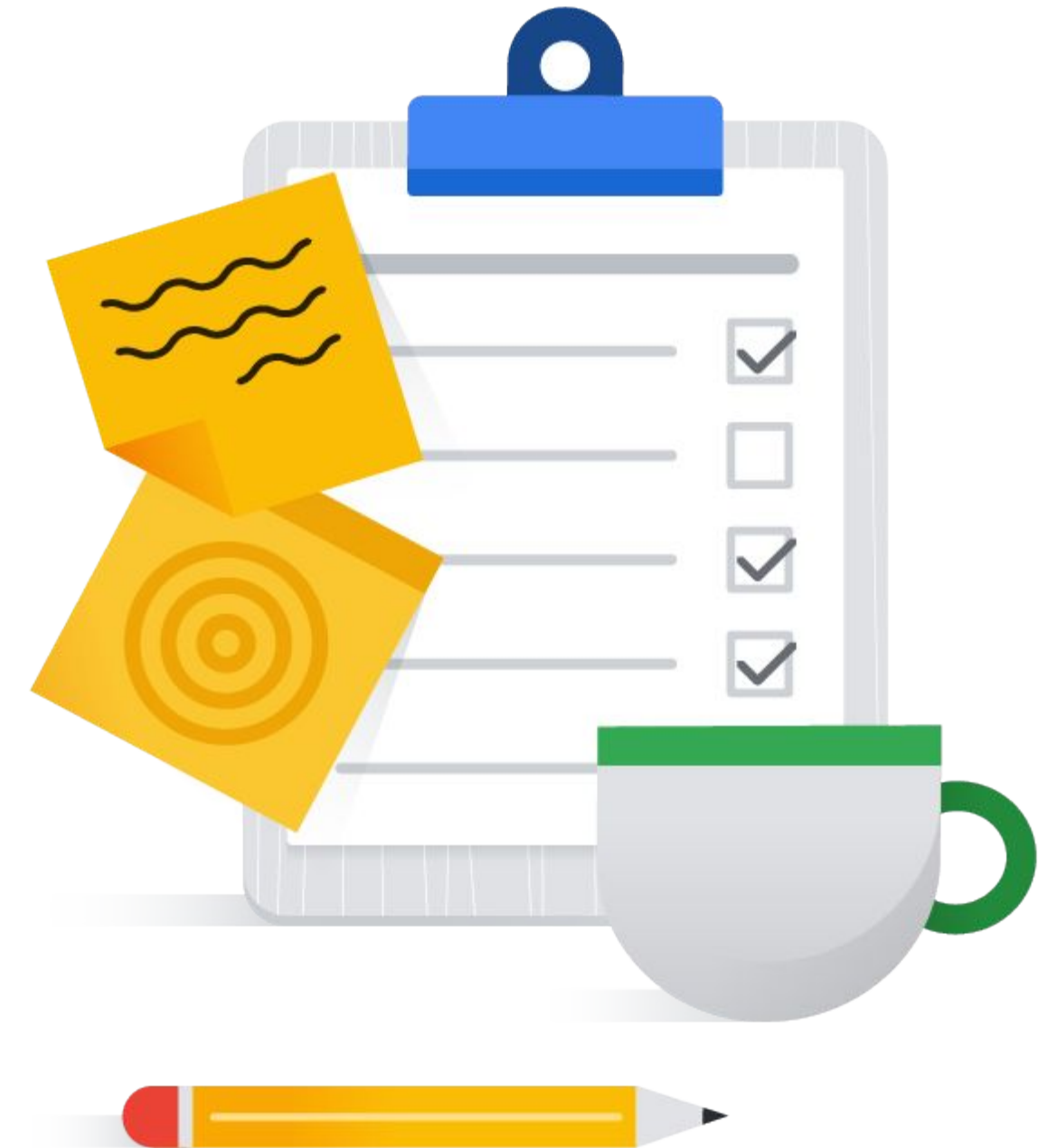## Security

Supports IAM, VPC-SC, and CMEK.

**5**

## Cost-effective
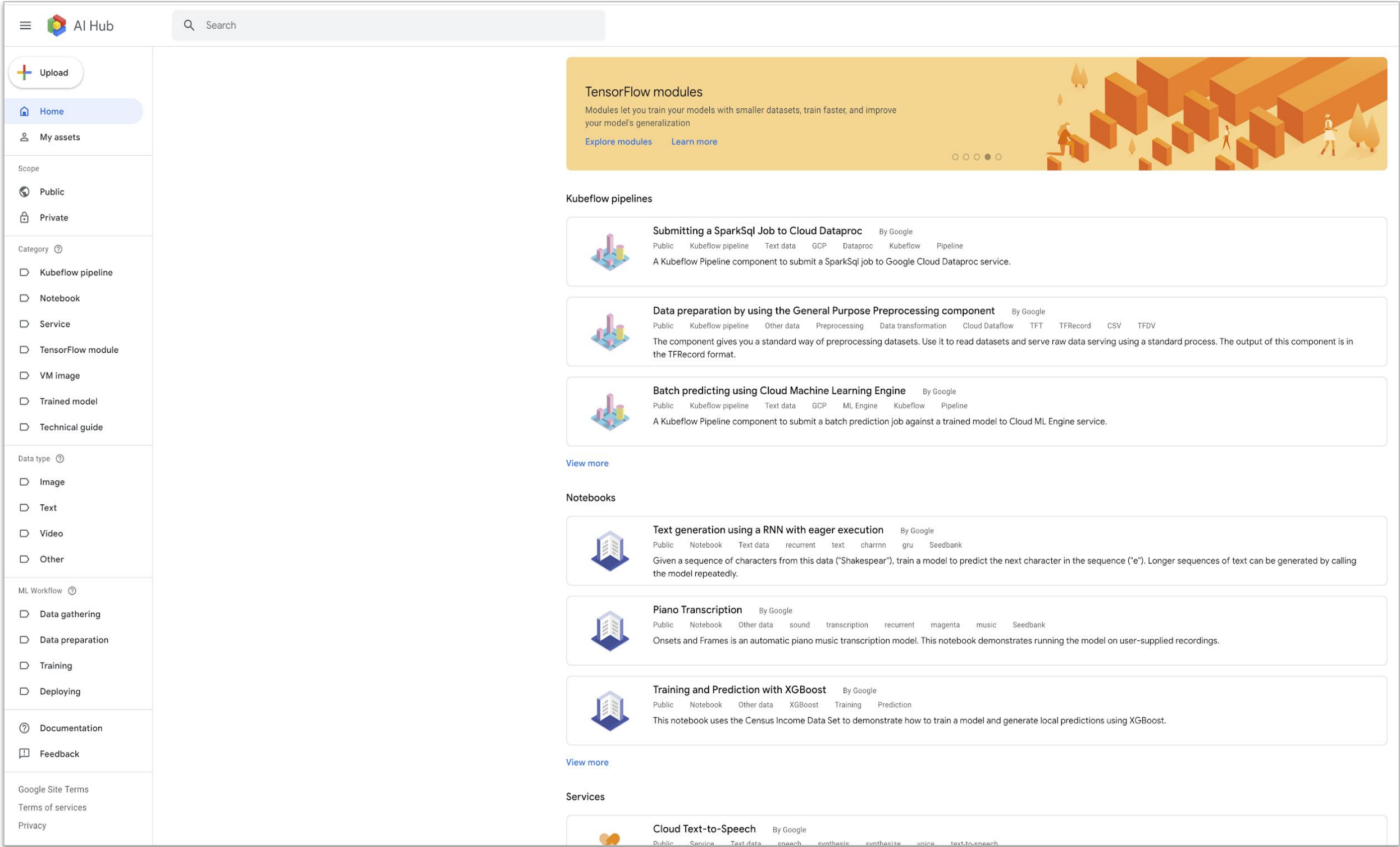
Only pay for the pipelines you run and the resources they use.

**6**

# Production ML Pipelines with Kubeflow

| | |
|---|---|
| 01 | Ways to do ML on Google Cloud |
| 02 | Vertex AI Pipelines |
| 03 | AI Hub |

# AI Hub is a repository for AI assets

Don't reinvent the wheel!
Find and deploy ML
pipelines.

# AI Hub stores various asset types

- Kubeflow pipelines and components

- Jupyter notebooks

- TensorFlow modules

- Trained models

- Services

- VM images

Google Cloud

# This is what a typical asset looks like

One-click deployment of ML pipelines via Kubeflow on Google Cloud as platform for AI, or on premise.
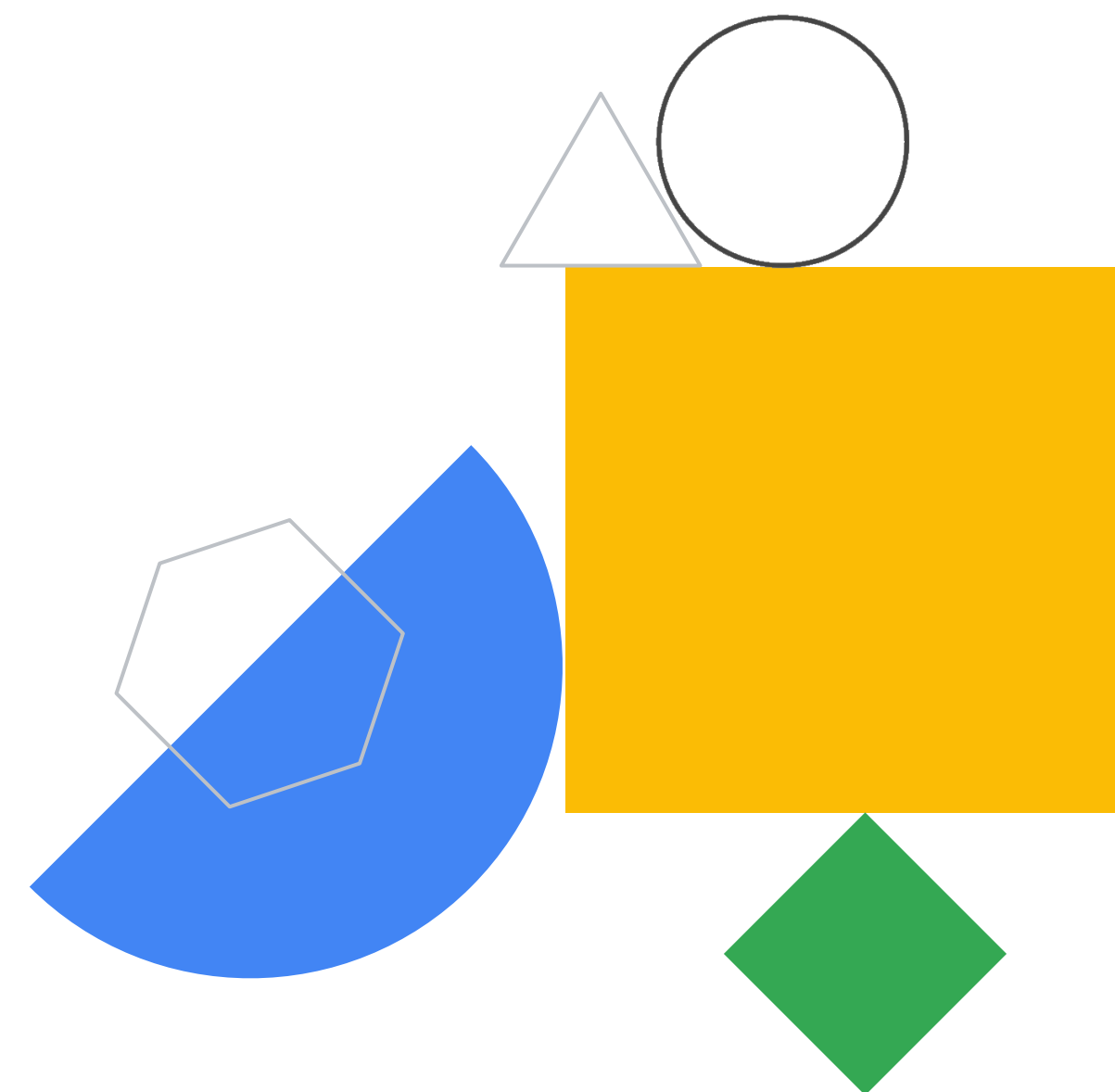
Google Cloud

# Assets on AI Hub are collected in two scopes: public assets and restricted assets

- Public scope are available to all AI Hub users.

- Restricted scope contains AI components that you have uploaded and assets that have been shared with you.

Google Cloud

# Lab Intro

Running Pipelines on Vertex AI

# Lab objectives

**01**    Set up the project environment

**02**    Inspect and configure pipeline code

**03**    Execute the AI pipeline

# Summary

- Use ML on Google Cloud using either:

  - Vertex AI (your model, your data)
  - AutoML (our models, your data)

- Use Vertex AI Pipelines to deploy end-to-end ML pipelines.

- Don't reinvent the wheel for your ML pipeline! Leverage pipelines on AI Hub.

Google Cloud