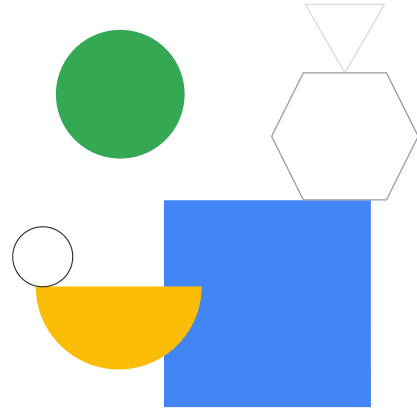# Production ML Pipelines with Kubeflow
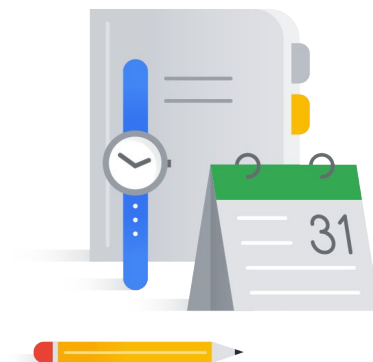
# Agenda

**Ways to do ML on Google Cloud**
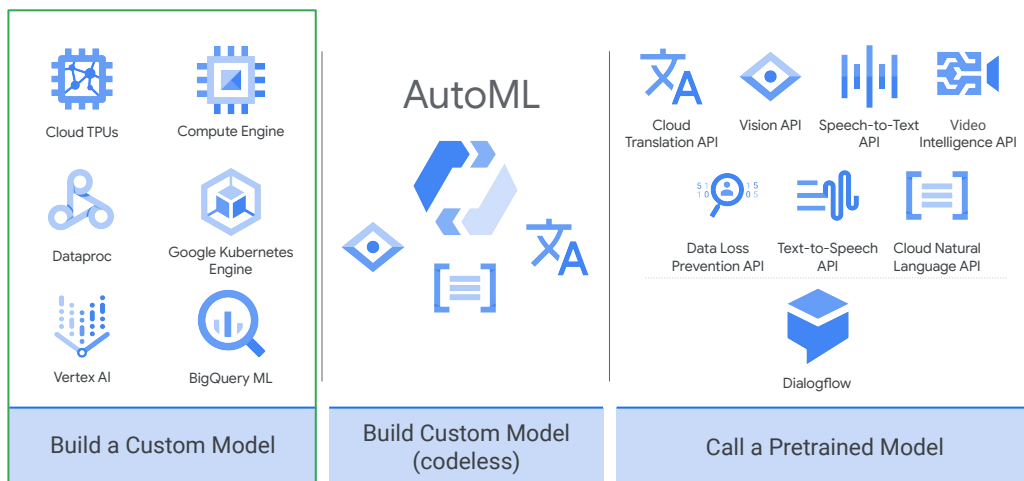
Kubeflow

AI Hub

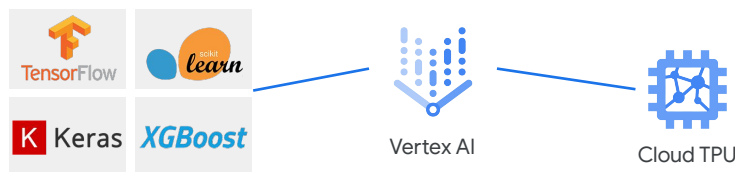Lab: Running ML Pipelines on Kubeflow

Google Cloud

In a previous module, we leveraged pre-trained ML APIs to process natural text. These are great options for seeing if your use case can just use a model that's already created and trained on Google's data. But, you may want a more tailored model trained on your own data. For that we will need a custom model. Let's talk about the different ways of building custom models.

# Create and deploy custom models with Kubeflow



| Build a Custom Model | Build Custom Model (codeless) | Call a Pretrained Model |

**Cloud TPUs** · **Compute Engine** · **Dataproc** · **Google Kubernetes Engine** · **Vertex AI** · **BigQuery ML**

**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**

As we covered previously, here are the three ways you can do machine learning on Google Cloud. We've already looked at pretrained models on the right. Now, we're going to visit the other side of the spectrum and build your own custom model and productionalize it on Google Cloud. There are a few ways of doing custom model development, training, and serving.

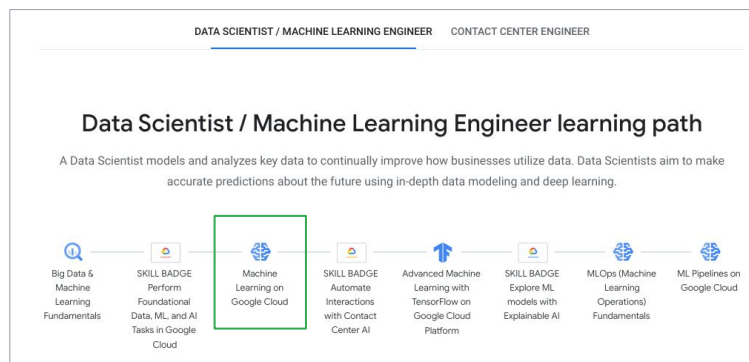What is Vertex AI exactly? It's a fully managed service for custom machine learning models, both training and serving predictions. It can scale from the experimentation stage all the way to production. You can also, using features of TensorFlow, include transformations on input data and perform hyperparameter tuning to choose the best model for your case. You can deploy your models to Vertex AI to serve predictions, which will autoscale to the demands of your clients.

Vertex AI also supports Kubeflow which is Google's open source framework for building ML pipelines -- and you'll have a lab on this later.

Essentially, Vertex AI is the engine behind doing machine learning at scale on Google Cloud. A data scientist can train and deploy production models from Notebooks with just a few commands.

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



DATA SCIENTIST / MACHINE LEARNING ENGINEER    CONTACT CENTER ENGINEER

## Data Scientist / Machine Learning Engineer learning path

A Data Scientist models and analyzes key data to continually improve how businesses utilize data. Data Scientists aim to make accurate predictions about the future using in-depth data modeling and deep learning.

Big Data & Machine Learning Fundamentals — SKILL BADGE Perform Foundational Data, ML, and AI Tasks in Google Cloud — Machine Learning on Google Cloud — SKILL BADGE Automate Interactions with Contact Center AI — Advanced Machine Learning with TensorFlow on Google Cloud Platform — SKILL BADGE Explore ML models with Explainable AI — MLOps (Machine Learning Operations) Fundamentals — ML Pipelines on Google Cloud

Google Cloud Training - Machine Learning and AI

Google Cloud

Since we're using Vertex AI and Kubeflow, we will often be thinking about using TensorFlow models. However, this isn't the course to dive into the details of TensorFlow. You can learn more about this in the Machine Learning on Google Cloud course, which is path of the Machine Learning and AI learning path for Data Scientists and Machine Learning Engineers linked here.

## Agenda

Ways to do ML on Google Cloud

Kubeflow

AI Hub

Lab: Running ML Pipelines on Kubeflow



Google Cloud

Where do Data Engineers come into the picture? Don't forget Data Engineers build data pipelines, and machine learning pipelines are no different. If we want to have a flexible pipeline for all stages of machine learning, Kubeflow is a great option.

## Perception: ML products are mostly about ML

Many people think that machine learning products are all about the code that ML scientists write locally on their machines. Does this code ensure the data going into it is clean? Can the code auto-scale to clients who want to use it for serving predictions? What if we have to re-train the model, does it go off-line at that point?

# Reality: ML requires lots of DevOps



Source: **Sculley et al.: Hidden Technical Debt in Machine Learning Systems**

The truth is, production machine learning systems are large, complicated, distributed systems. There's a lot of DevOps involved for things like monitoring, and process management tools. Google started building Kubeflow to tackle these DevOps challenges using Kubernetes and containers.

# Kubeflow provides a platform for building ML products

- Leverage containers and Kubernetes to solve the challenges of building ML products.

- Kubeflow = Cloud Native, multi-cloud solution for ML.

- Kubeflow provides a platform for composable, portable and scalable ML pipelines.

- If you have a Kubernetes conformant cluster, you can run Kubeflow.

Google Cloud

# Kubernetes is a great platform for ML

- Containers
- Scaling built in
- Unified architecture
- Easy to integrate building blocks
  - ML APIs
  - Dataflow
- Lots of options for CI/CD
- Portability
  - Dev, On-Prem, Multi-cloud: same stack

Google Cloud

Before we dive into the specifics of Kubeflow Pipelines, I should provide additional context.

- Kubeflow Pipelines is a part of the open source project Kubeflow.
- Kubeflow is a platform that provides the tools and scalable services required to develop and deploy ML workloads, all the way from distributed training, to scalable serving, to Notebooks w/ JupyterHub and workflow orchestration and much more.
- Kubeflow services are built on top on Kubernetes. Kubernetes provides scalability and hybrid protability. You can run Kubeflow anywhere you can run a Kubernetes cluster, and thus applications built on Kubeflow are portable across clouds and on-premise environments. On Google Cloud, you can easily deploy Kubeflow on Google Kubernetes Engine.

## The capabilities provided by Kubeflow pipelines can largely be put into three buckets

**ML Workflow Orchestration**

**Share, Re-use & Compose**

**Rapid Reliable Experimentation**

One option to help manage the overhead of productionalizing ML pipelines is to use Kubeflow

The capabilities provided by Kubeflow pipelines can largely be put into three buckets:

- ML Workflow Orchestration
- Share, Re-use & Compose
- Rapid Reliable Experimentation

You can think of the benefits as similar to those of Cloud Composer but better tailored for ML workloads. Let's see what a pipeline looks like.

# What constitutes a Kubeflow pipeline

**Containerized implementations of ML Tasks**

- Containers provide portability, repeatability and encapsulation.
- A task can be single node or *distributed*
- A containerized task can invoke other services like Vertex AI, Dataflow, or Dataproc.

**Specification of the sequence of steps**

- Specified via Python SDK.

**Input Parameters**

- A "Job" = Pipeline invoked w/ specific parameters.

Let me provide a peek under the hood.

## Visual depiction of pipeline topology

To make things more concrete let's look at a screenshot of an illustrative workflow that was run on Kubeflow Pipelines. This is just an illustrative workflow and users can author and run many different kinds of workflow topologies with different code and tools in the various steps of the workflow.

For each workflow that is run on Kubeflow Pipelines, you get a rich visual depiction of the topology so that you know what was executed as part of the workflow.

1. In this workflow we start with a data preprocessing and validation step
2. Followed by feature engineering
3. Following by a fork where we traing many different kinds of models
4. The models that are trained are then analyzed and compared on a test data set
5. Finally, if an improved model is produced, it is deployed to a serving endpoint

# Rich visualization of metrics

For each step of the workflow, you have rich ML specific information at your fingertips. Just click on a step and visualize relevant metrics produced by that step, such as an ROC curve for example. If you did model training and produced the rich metadata that can be visualized with TensorBoard, that is just a click away.

# View all configs, inputs and outputs



For each step of the workflow you can see the precise configuration parameters, inputs and outputs. Thus, for a model trained with Kubeflow Pipelines you never have to wonder, how exactly did I create this model?

Here you can quickly see how long the model training took, where the trained model is, and what data was used for training and evaluation.

# Author pipelines with an intuitive Python SDK

You can define the ML workflow using Kubeflow's Python SDK. By defining the workflow we mean specifying each step's inputs and outputs and how the various steps are connected. The topology of the workflow is implicitly defined by connecting the outputs of an upstream step to the inputs of a downstream step. You can also define looping constructs as well as conditional steps.

___

# Package and share pipelines as zip files

- Upload and execute pipelines via UI (in addition to API/SDK).

- Pipeline steps can be authored as reusable components.



Google Cloud

Another nice Kubeflow feature is the ability to package pipeline components. This adds an element of portability since you can then move your ML pipelines even between cloud providers.

Kubeflow pipelines separate the work for different parts of the pipeline to enable people to specialize. For example, a ML engineer can focus on feature engineering and the other parts of creating the model such as hyperparameter tuning. The ML engineer's solutions can then be bundled up and used by a data engineer as part of a data engineering solution. The solution can then appear as a service used by a data analyst to derive business insights.

# Rapid, reliable, experimentation

- Every run logged with all config params, inputs, outputs and metrics.
- Easily search and find old runs.
- Clone and re-run or modify.

Clone

Edit

**Rapidly iterate on ideas**

Submit

# View all current and past runs in one place

# Easy comparison and analysis of runs

# Easy comparison and analysis of runs

Kubeflow makes it easy to run a number of ML experiments at the same time. For example, if you're doing hyperparameter optimization, you can easily deploy a number of different training instances with different hyperparameter sets. Kubeflow's run overview makes it easy to hone in on the techniques or parameters generating the best results. You can quickly identify what worked and what did not work.
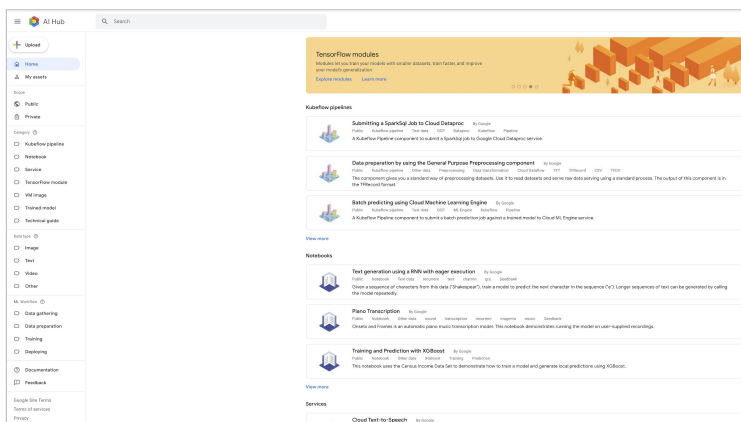
# Agenda

Ways to do ML on Google Cloud

Kubeflow

AI Hub

Lab: Running ML Pipelines on Kubeflow

Google Cloud

We mentioned that Kubeflow pipelines can be packaged and shared with other users. This leads us to a discussion of AI Hub.

# AI Hub is a repository for AI assets

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



Google Cloud

AI Hub is a repository for ML components. Don't reinvent the wheel! Avoid building some component when someone else has already built it, and most likely, has already optimized it. You can find and deploy not just containerized applications for machine learning, but full ML pipelines on AI hub.

# AI Hub stores various asset types

- Kubeflow pipelines and components
- Jupyter notebooks
- TensorFlow modules
- Trained models
- Services
- VM images

Google Cloud

What asset types can we find on AI Hub? Among the assets stored on AI Hub are entire Kubeflow pipelines, Jupyter notebooks, TensorFlow modules, fully trained models, services and VM images.

# This is what a typical asset looks like



Here you see what a typical asset looks like. You can see information about the pipeline, such as inputs and outputs, and download options.

___

## 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.

The assets on AI Hub are collected into two scopes: public assets and restricted assets. Public assets are available to all AI Hub users. Restricted scope assets contains AI components you have uploaded and those that have been shared with you. For example, you could have assets only available to people within your organization or teams.
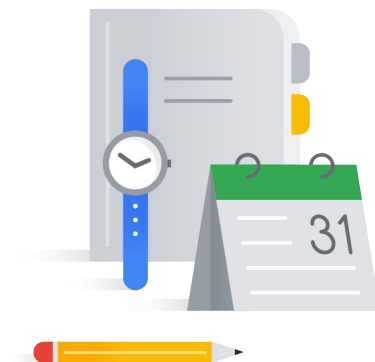
# Agenda

Ways to do ML on Google Cloud

Kubeflow

AI Hub

Lab: Running ML Pipelines on Kubeflow

# Lab Intro

## Running ML Pipelines on Kubeflow

In this lab you learn how to install and use Kubeflow Pipelines. The objectives of the lab are for you to:

- Create a Kubernetes cluster and configure AI Platform pipelines.
- Launch the pipelines dashboard.
- Create and run an experiment from an example end-to-end ML Pipeline.
- Examine and verify the output of each step.
- Inspect the pipeline graph, various metrics, logs, charts and parameters.

# Module summary

- Use ML on Google Cloud using either:
  - Vertex AI (your model, your data)
  - AutoML (our models, your data)
  - Perception API (our models, our data)

- Use Kubeflow to deploy end-to-end
  ML pipelines.

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

Google Cloud