# S91030 - Hybrid Machine Learning with the Kubeflow Pipelines and RAPIDS
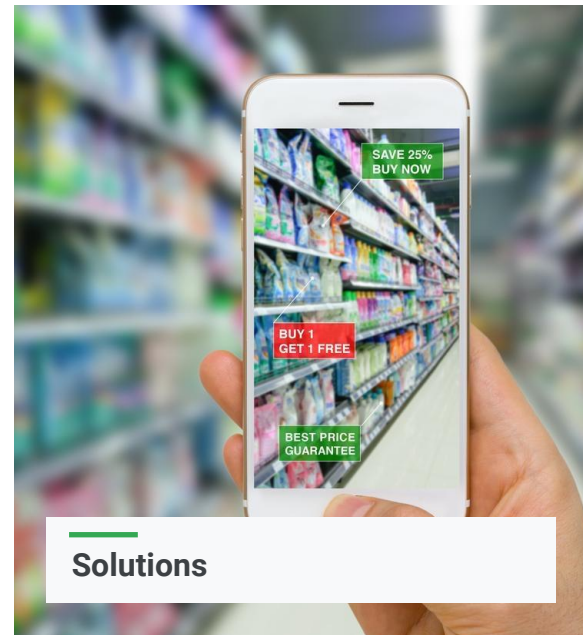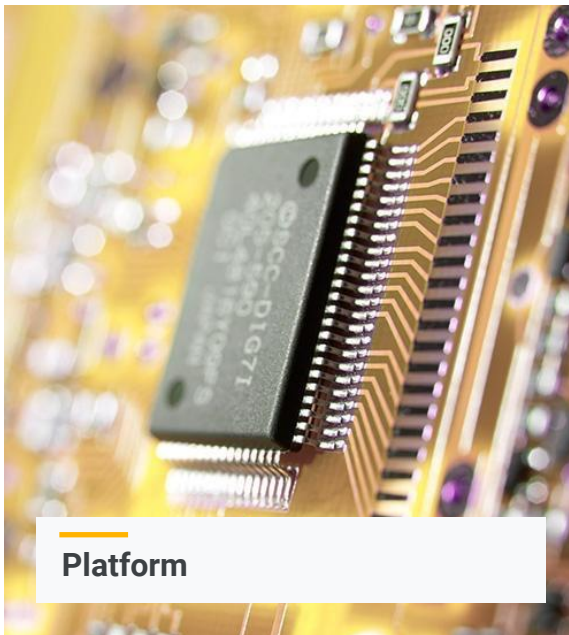
**Sina Chavoshi**
Technical Program Manager

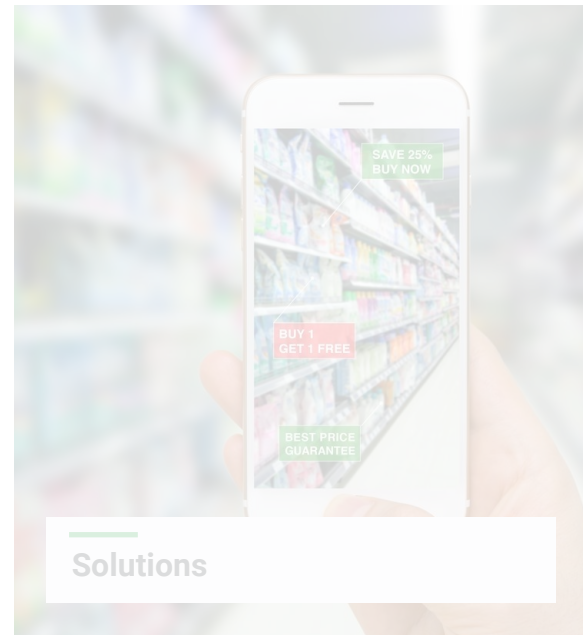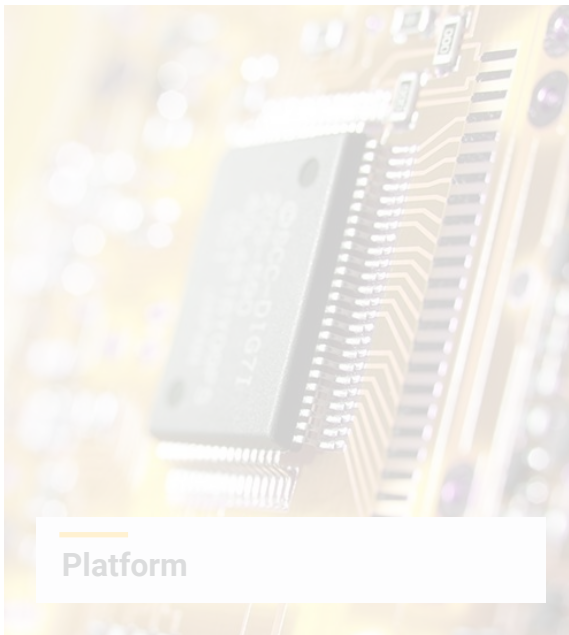# Cloud AI Strategy:
## The right approach for the right problem



**Building blocks**



**Platform**



**Solutions**

Google Cloud

# Cloud AI Strategy:

## The right approach for the right problem



**Building blocks**



Platform



Solutions

Google Cloud

# Building Blocks

## Sight

**Cloud Vision API**
Image recognition and classification.

**Cloud Video Intelligence API**
Scene-level video annotation.

**AutoML Vision** BETA
Custom image classification models.

## Language

**Cloud Translation API**
Language detection and translation.

**Cloud Natural Language API**
Text parsing and analysis.

**AutoML Translation** BETA
Custom domain-specific translation.

**AutoML Natural Language** BETA
Custom text classification models.

## Conversation

**Dialogflow Enterprise Edition**
Build conversational interfaces.
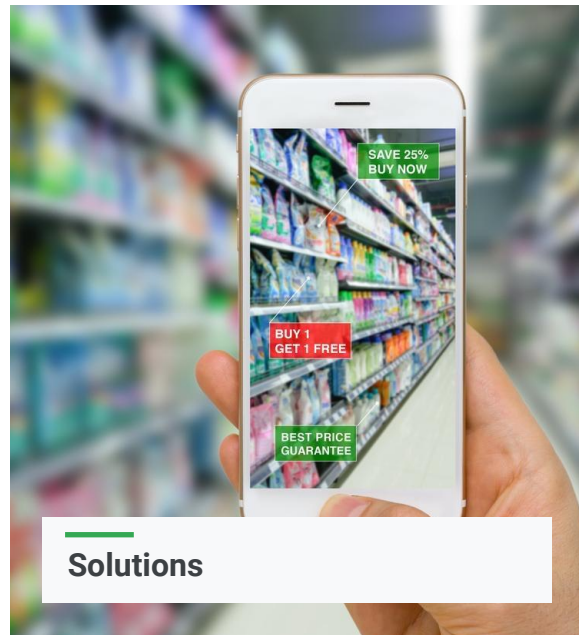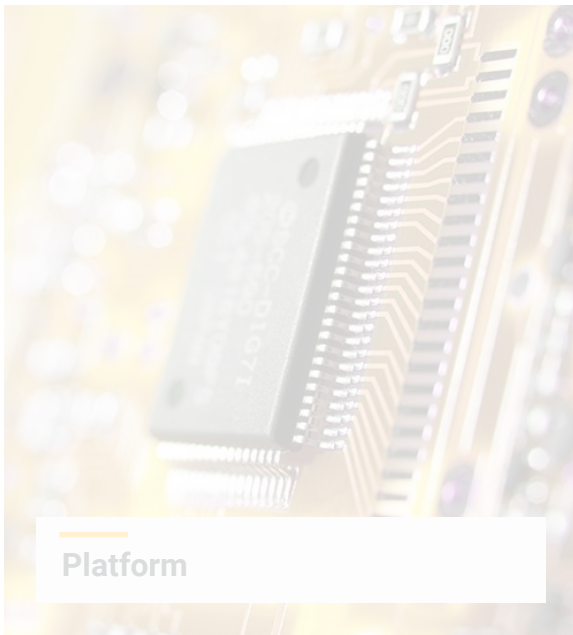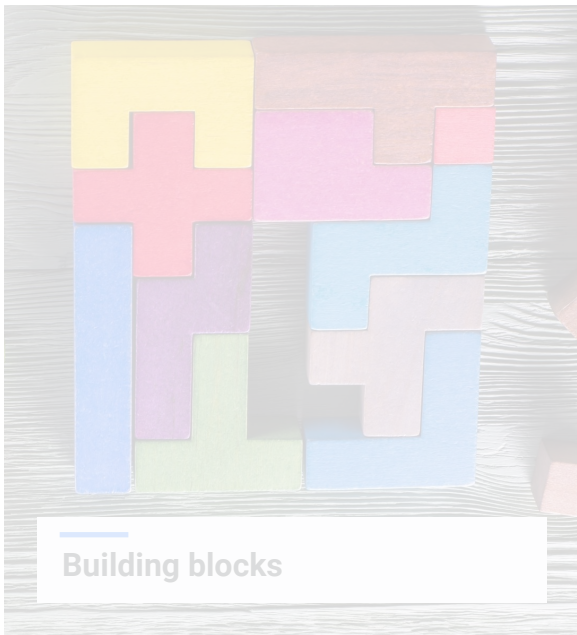
**Cloud Text-to-Speech API**
Convert text to speech.

**Cloud Speech-to-Text API**
Convert speech to text.

Google Cloud

# Cloud AI Strategy:
## The right approach for the right problem



Building blocks



Platform



Solutions

Google Cloud

# Solutions / Contact Center



Phone

Customer

Chat

**Google Cloud Contact Center AI**

Contact Center Provider

Contact Center Interface

Virtual Agent

Backend Fulfillment

Knowledge Base (PDF/HTML)

Agent Assist

Virtual Agent

Agent

Google Cloud

# Cloud AI Strategy:
## The right approach for the right problem



Building blocks

Platform

Solutions

Google Cloud

# Cloud AI Platform

Building & deploying real-life ML applications is **hard** and **costly** because of **lack of tooling** that covers **end-to-end ML** development & deployment.

Google Cloud

# In addition to the actual ML...



ML
Code

Google Cloud

# You have to worry about so much more.



Configuration

Data Collection

Data Verification

Monitoring

Serving Infrastructure

ML Code

Analysis Tools

Feature Extraction

Process Management Tools

Machine Resource Management

TensorFlow

Google Cloud

# AI **problems** today

## Problems

### Deployment
Brittle, opinionated infrastructure that is hard to productionize and breaks between cloud and on-prem

### Talent
Machine Learning expertise is scarce

### Collaboration
Difficult to find, leverage existing solutions

## Solutions

**01** Kubeflow

**02** Reusable pipelines

**03** Google Cloud AI Hub

Google Cloud

# 01: **Kubeflow**

**ML microservices**

## Scalable ML services on Kubernetes

**Easy to get started**
- Out-of-box support for top frameworks
  - pytorch, caffe, tf and xgboost
- Kubernetes manages dependencies, resources

**Swappable & scalable**
- Library of ML services
- GPU support
- Massive scale

**Meet customer where they are**
- GCP
- On-prem with Cisco

Cloud

On-prem

Training        Predict

Training        Predict

**kubernetes**

**Google** Cloud

**NVIDIA**

# RAPIDS

Product Overview

# THE BIG PROBLEM IN DATA SCIENCE

# BENCHMARKS

### cuIO/cuDF —
### Load and Data Preparation



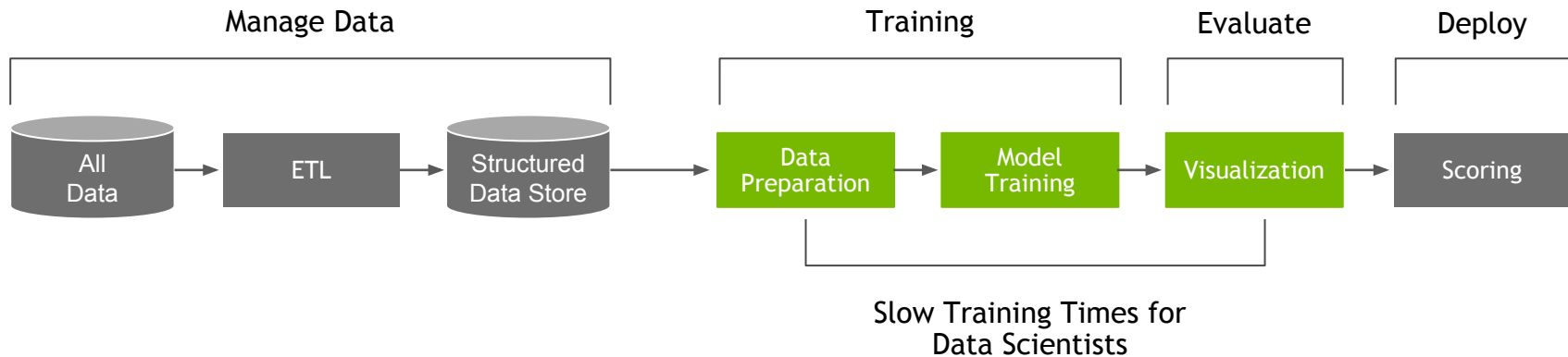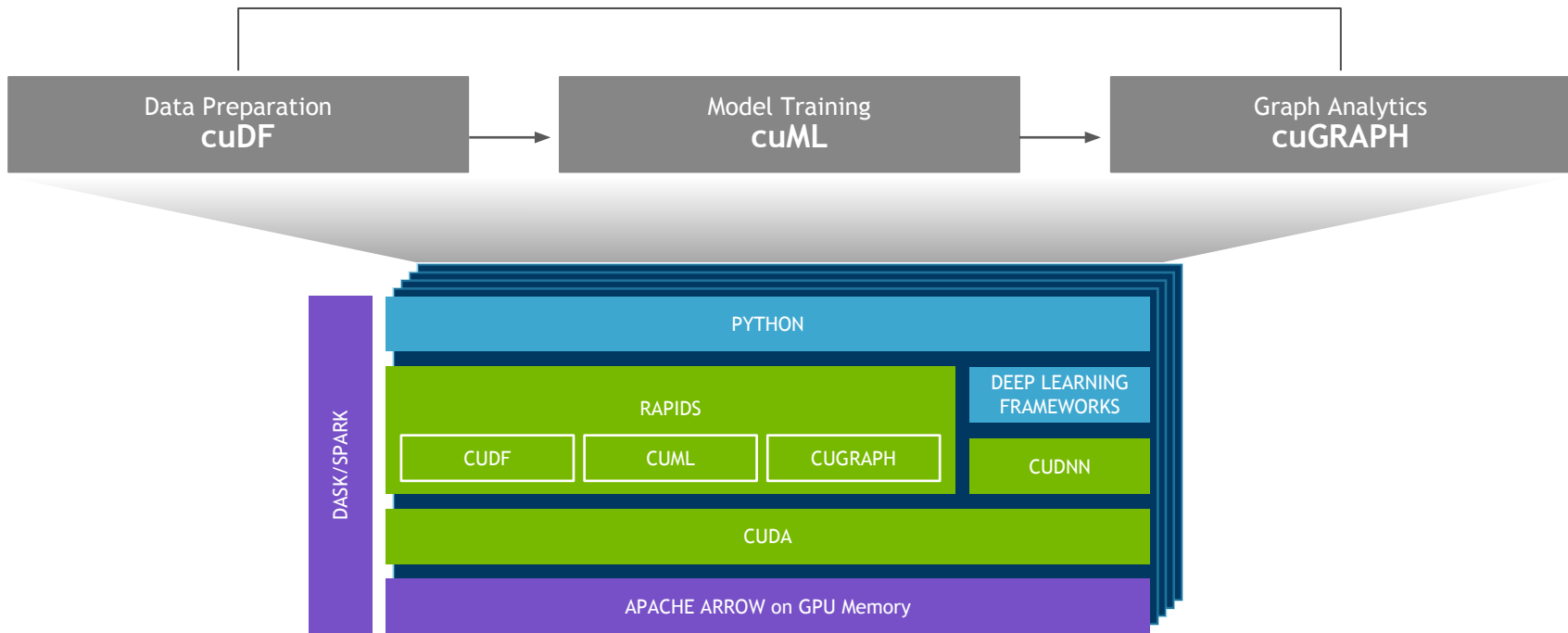| | |
|---|---|
| 20 CPU Nodes | 2741 |
| 30 CPU Nodes | 1675 |
| 50 CPU Nodes | 715 |
| 100 CPU Nodes | 379 |
| DGX-2 | 42 |
| 5x DGX-1 | 19 |

### cuML — XGBoost

| | |
|---|---|
| 20 CPU Nodes | 2290 |
| 30 CPU Nodes | 1956 |
| 50 CPU Nodes | 1999 |
| 100 CPU Nodes | 1948 |
| DGX-2 | 169 |
| 5x DGX-1 | 157 |

**Time in seconds — Shorter is better**

### End-to-End

| | |
|---|---|
| 20 CPU Nodes | |
| 30 CPU Nodes | |
| 50 CPU Nodes | |
| 100 CPU Nodes | |
| DGX-2 | |
| 5x DGX-1 | |

■ cuIO / cuDF (Load and Data Preparation)    Data Conversion
XGBoost

**Benchmark**

200GB CSV dataset; Data preparation includes joins, variable transformations.

**CPU Cluster Configuration**

CPU nodes (61 GiB of memory, 8 vCPUs, 64-bit platform), Apache Spark

**DGX Cluster Configuration**

5x DGX-1 on InfiniBand network

Google Cloud

# AI Hub & Pipelines: Fast & simple adoption of AI



The Flywheel of AI Adoption

**1. Search & Discover**
Find best-of-breed solutions on the AI Hub which leverage Cloud AI solutions

**2. Deploy**
Quick 1-click implementation of ML pipelines onto Google Cloud Platform .

**3. Customize**
Experiment and adjustment out-of-the-box pipelines to custom use cases.

**4. Run in production**
Deploy customized pipelines in production.

**5. Publish**
Upload & share pipelines running best within your org or publicly.

Network effect

Google Cloud

# 02: **Reusable Pipelines**

**Enable developers to build custom ML applications by easily "stitching" and connecting various components.**

- Reuse instead of reimplement or reinvent
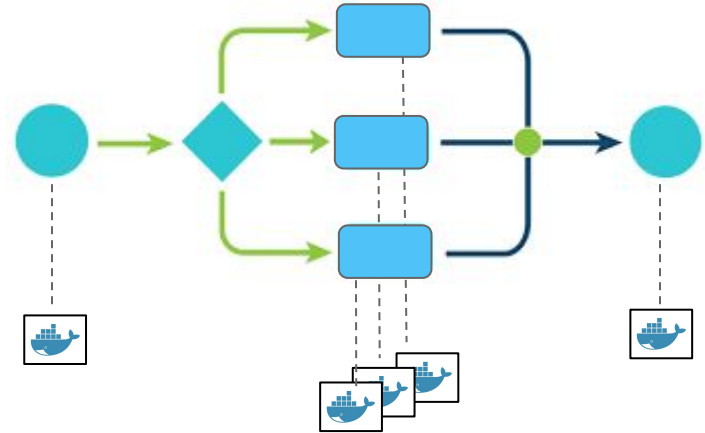- Discover, learn and replicate successful pipelines

Google Cloud

# 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

- Specification of the sequence of steps
  - Specified via Python SDK

- Input Parameters
  - A "Job" = Pipeline invoked w/ specific parameters



Google Cloud

# 03: AI Hub at a glance

**1** **All AI content in one place**
**Quick discovery** of **plug & play** AI pipelines & other content built by teams across Google and by partners and customers.

**2** **Fast & simple implementation of AI on GCP**
One-click deployment of AI pipelines via Kubeflow on **GCP as the go-to platform for AI** + hybrid & on premise.

**3** **Enterprise-grade internal & external sharing**
**Foster reuse** by sharing deployable AI pipelines & other content privately within organizations & publicly.



## Google Cloud

# Mission

The one place for everything
AI, from experimentation
to production.

Google Cloud

# Public and private AI Hub

🔓 **Public content**

🔒 **+ Private content**

**By Google**

Unique AI assets by Google

**By partners**

Created, shared & monetized by anyone

**By customers**

Content shared securely within and with other organizations

AutoML, TPUs, kaggle Cloud AI Platform, etc.

Research at Google

DeepMind

Google Cloud

# Kubeflow Pipelines enable



**Workflow orchestration**



**Rapid reliable experimentation**



**Share, re-use & compose**

Google Cloud

# Demo

Google Cloud

← ✔ ra-run-2b7

Clone    Archive

Graph     Run output     Config



Tensorflow-Wide-and-Deep ✔

PreProcess ✔ → FeatureTransforms ✔ → XGBoost-GBT ✔ → ModelValidation ✔ → PushToServing ✔

Tensorflow-CNN ✔

**Visual depiction of pipeline topology**

ⓘ Pipeline runtime graph

# Experiments

**All experiments**    All runs

Filter experiments

| | Experiment name | Last 5 runs | Created on ↑ | Created by |
|---|---|---|---|---|
| ▶ | tfma-experiment | ↻ | 6:17 PM, Aug 24, 2018 | John Doe |
| ▶ | xgboost-train | ✓ ✓ ✓ | 6:17 PM, Aug 24, 2018 | John Doe |
| ▶ | promo-email | ✓ ✓ ✓ ✓ ! | 6:17 PM, Aug 24, 2018 | Walter Fisher |
| ▶ | data-prep | ✓ | 6:17 PM, Aug 24, 2018 | Walter Fisher |
| ▶ | tf-preprocessing | ✓ ✓ ✓ | 6:17 PM, Aug 24, 2018 | John Doe |
| ▶ | tf-training | ✓ ✓ ✓ ✓ ! | 6:17 PM, Aug 24, 2018 | Walter Fisher |
| ▶ | tf-serving | ✓ ✓ ✓ | 6:17 PM, Aug 24, 2018 | Walter Fisher |

Rows per page: 10 ▾     1–10 of 241    ‹ ›

**View all current and historical runs, grouped as "Experiments"**

← ✓ ccard-recommender-run1

Clone     Archive

Graph     Run output     Config

✕  Train

Artifacts     Logs     Config

ROC curve

TPR

1.0
0.8
0.6
0.4
0.2
0

0  0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0

FPR

Tensorboard

⧉ Open Tensorboard

PreProcess ✓  →  FeatureTransforms ✓  →  Train ✓  →  Mo...

**Rich visualizations of metrics**

ⓘ Pipeline runtime graph

Clone    Archive

fig



X    Train

Artifacts    Logs    Config

ROC curve

FeatureTransforms    Train    Mc

**Clone an existing pipeline**

Tensorboard

← ✅ **Simple XGBoost Classifier**

**Graph**     **Config**

## Run details

| | |
|---|---|
| **Status** | Succeeded |
| **Description** | |
| **Created at** | 11/25/2018, 12:56:44 PM |
| **Started at** | 11/25/2018, 12:56:44 PM |
| **Finished at** | 11/25/2018, 1:16:37 PM |
| **Duration** | 0:19:53 |

## Run parameters

| | |
|---|---|
| **output** | gs://mlpipelines |
| **project** | foo2thebar |
| **region** | us-central1 |
| **train-data** | gs://ml-pipeline-playground/sfpd/train.csv |
| **eval-data** | gs://ml-pipeline-playground/sfpd/eval.csv |
| **schema** | gs://ml-pipeline-playground/sfpd/schema.json |
| **target** | resolution |
| | |
| **true-label** | ACTION |

**Access to all config params, inputs and outputs for each run**

Run details

Pipeline *
xgboost training - confusion matrix                          Choose

Run name *
product-recommender-model

Description (optional)
Train XBG model for product recommendation application.

Run parameters

Specify parameters required by the pipeline

output

project

region
us-central1

train-data
gs://ml-pipeline-playground/sfpd/train.csv

eval-data
gs://ml-pipeline-playground/sfpd/eval.csv

schema
gs://ml-pipeline-playground/sfpd/schema.json

target
resolution

rounds
200

workers
2

true-label
ACTION

image-classifier

Edit    Archive

Fastest run time | Slowest run time
1m 59s | 3m 20s
View run | View run

## All runs

Start new run    Start recurring run    Compare runs    Stop    Archive    Metrics

Filter runs

| | Runs | Status | Duration | Pipeline | Recurring run config. | Start time ↑ | rmse | eta |
|---|---|---|---|---|---|---|---|---|
| ☐ | ccard-recommender-run3 | ✅ | 1m 59s | linear-classifier | | 9:32 AM, Aug 26, 2018 | 0.88 | 0.92 |
| ☐ | ccard-recommender-run2-clone(2) | ✅ | 2m 12s | linear-classifier | | 11:42 AM, Aug 25, 2018 | 0.72 | 0.86 |
| ☐ | ccard-recommender-run2-clone(1) | ✅ | 2m 44s | linear-classifier | | 10:48 AM, Aug 25, 2018 | 0.74 | 0.84 |
| ☐ | ccard-recommender-run2 | ✅ | 2m 18s | linear-classifier | | 10:22 PM, Aug 25, 2018 | 0.82 | 0.76 |
| ☐ | ccard-recommender-run1-clone(1) | ✅ | 2m 20s | linear-classifier | | 10:10 AM, Aug 25, 2018 | 0.80 | 0.84 |
| ☐ | ccard-recommender-run1 | ✅ | 3m 20s | linear-classifier | | 6:17 PM, Aug 24, 2018 | 0.72 | 0.76 |

Rows per page:    10 ▾    1–10 of 241    ‹    ›

**Easy comparison of Runs**

← image-classifier

Edit    Archive

Fastest run time
**1m 59s**
View run

Slowest run time
**3m 20s**
View run

## All runs

Start new run   Start recurring run   Compare runs   Stop   Archive   Metrics

Filter runs

| | Runs | Status | Duration | Pipeline | Recurring run config. | Start time ↑ | rmse | eta |
|---|---|---|---|---|---|---|---|---|
| ☑ | ccard-recommender-run3 | ✓ | 1m 59s | linear-classifier | | 9:32 AM, Aug 26, 2018 | 0.88 | 0.92 |
| ☑ | ccard-recommender-run2-clone(2) | ✓ | 2m 12s | linear-classifier | | 11:42 AM, Aug 25, 2018 | 0.72 | 0.86 |
| ☑ | ccard-recommender-run2-clone(1) | ✓ | 2m 44s | linear-classifier | | 10:48 AM, Aug 25, 2018 | 0.74 | 0.84 |
| ☐ | ccard-recommender-run2 | ✓ | 2m 18s | linear-classifier | | 10:22 PM, Aug 25, 2018 | 0.82 | 0.76 |
| ☐ | ccard-recommender-run1-clone(1) | ✓ | 2m 20s | linear-classifier | | 10:10 AM, Aug 25, 2018 | 0.80 | 0.84 |
| ☐ | ccard-recommender-run1 | ✓ | 3m 20s | linear-classifier | | 6:17 PM, Aug 24, 2018 | 0.72 | 0.76 |

Rows per page:    10 ▾    1–10 of 241    ‹  ›
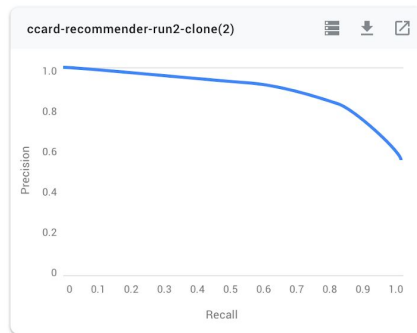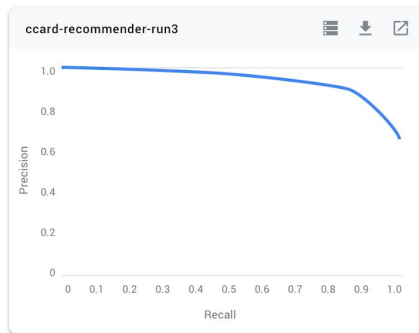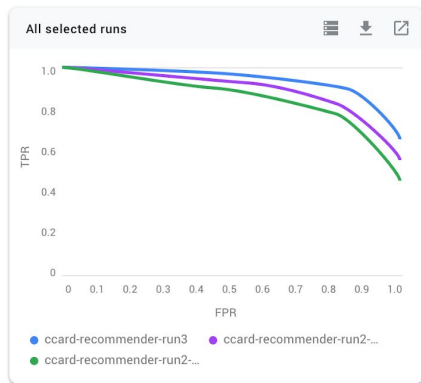
**Easy comparison of Runs**

← Compare runs

## Run overview

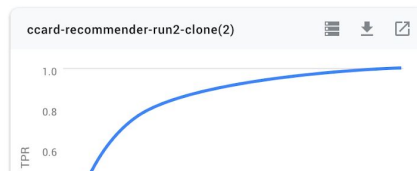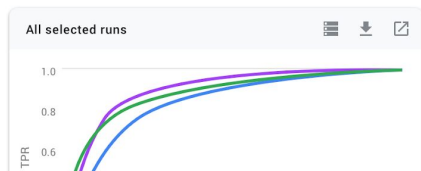| Show | Run name | Status | Pipeline | Duration |
|------|----------|--------|----------|----------|
| ☑ | ccard-recommender-run3 | ✔ | linear-classifier | 3m 20s |
| ☑ | ccard-recommender-run2-clone(2) | ✔ | linear-classifier | 3m 20s |
| ☑ | ccard-recommender-run2-clone(1) | ✔ | linear-classifier | 3m 20s |

▸ Parameters

▸ Metrics

▾ Precision Recall



▾ ROC curve

# That's a wrap.

Google Cloud