# Containerized or Serverless Learning

**Comparing the Prediction Latency of Machine Learning Model Deployments on Containerized and Serverless Environments**

Diya Ramasamy | Iris Ma | Krithika Balasubramanyam
Raj Shreyas Penukonda | Yuyu Lai

# System Purpose and Functionality

Deployed three different machine learning models on containerized and serverless environments as APIs
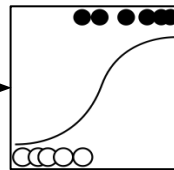
**Logistic Regression**
(Scikit-Learn)
➔ Input: Tabular Data
➔ Output: Iris Species

**Computer Vision Model**
(TensorFlow)
➔ Input: Grayscale Image
➔ Output: Clothing Type

**NLP Model** (TensorFlow)
➔ Input: Wine Review
➔ Output: Type of Wine

# Design

## Pipeline



Google Cloud Bucket

Google Cloud Functions
-Serverless API-

Vertex AI
-Containerized API-

POSTMAN

## Metrics

| EVENT | TIME | EVENT | TIME |
|---|---|---|---|
| Prepare | 6.81 ms | Prepare | 7.42 ms |
| Socket Initialization | 2.49 ms | Socket Initialization | 1.38 ms |
| DNS Lookup | 6.69 ms | DNS Lookup | Cache |
| TCP Handshake | 63.16 ms | TCP Handshake | Cache |
| SSL Handshake | 89.81 ms | SSL Handshake | Cache |
| Transfer Start | 257.5 ms | Transfer Start | 150.72 ms |
| Download | 16.16 ms | Download | 6.54 ms |
| Process | 9.98 ms | Process | 5.13 ms |
| Total | 452.57 ms | Total | 171.17 ms |

| Batch Sizes | Transfer Time (#1) | Transfer Time (#2) | Transfer Time (#3) |
|---|---|---|---|
| 1<br>2<br>5<br>10<br>50<br>100 | … | … | … |

# Results and Challenges



Logistic Regression: Vertex AI vs Cloud Functions Transfer Time

Computer Vision: Vertex AI vs Cloud Functions Transfer Time

Natural Language Processing: Vertex AI vs Cloud Functions Transfer Time

- Vertex AI - Trial 1 (Not Cached)
- Vertex AI - Trial 2 (Cached)
- Vertex AI - Trial 3 (Cached)
- Cloud Functions - Trial 1 (Not Cached)
- Cloud Functions - Trial 2 (Cached)
- Cloud Functions - Trial 3 (Cached)

GCF outperforms Vertex AI; dynamic resource allocation allows for consistently low latency

| Challenge | Resolution |
|---|---|
| Out of memory issues for large models | Trained another model from scratch with Tensorflow |
| Framework version conflicts | Trained models with the exact version available on GCP |
| Authentication issues when calling endpoints | Regenerated tokens every hour |