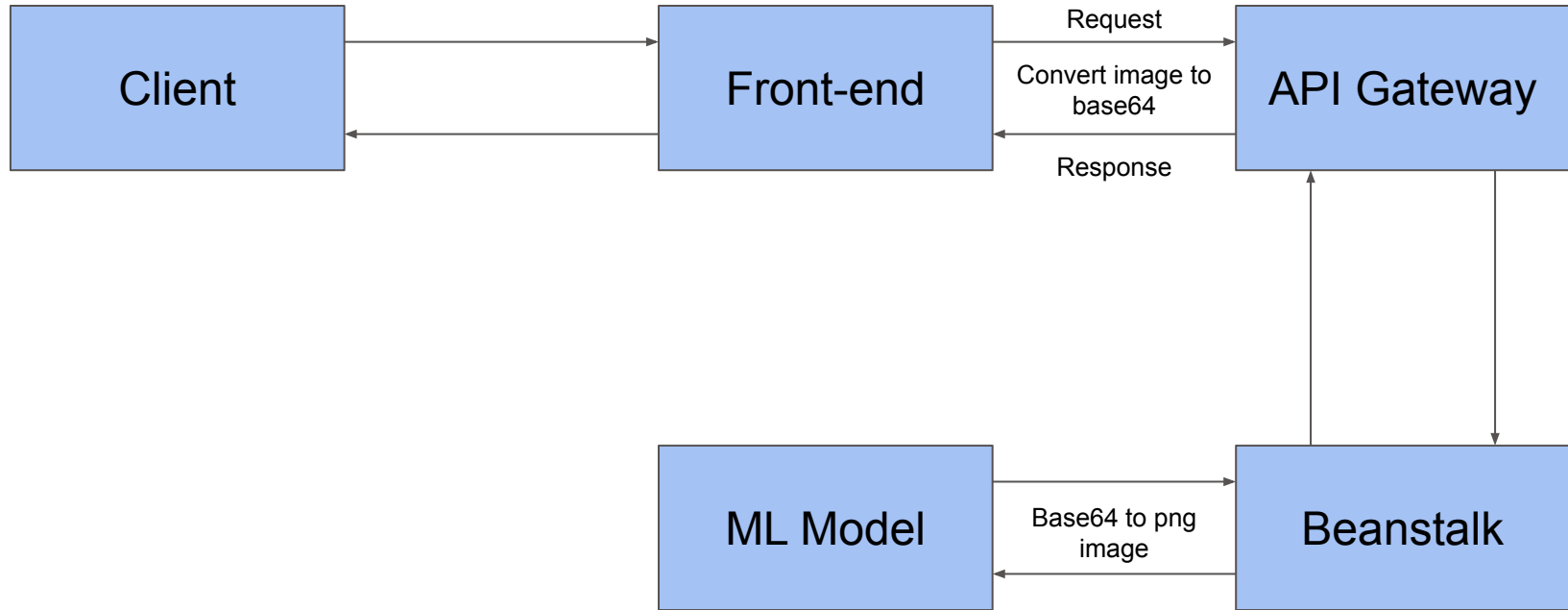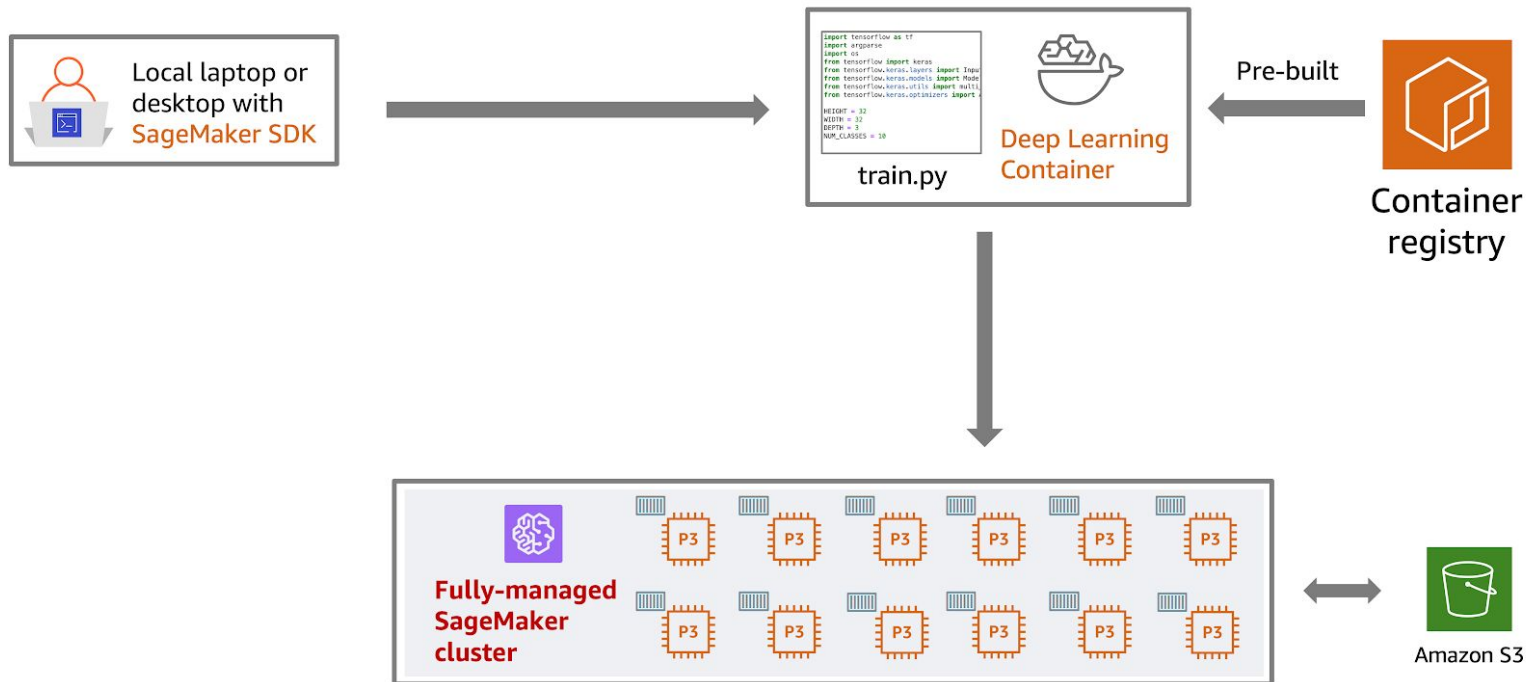# Distributed Machine Learning

# Idea

Perform distributed training with TensorFlow and Horovod on Amazon Sagemaker for image classification of fruit dataset from Kaggle
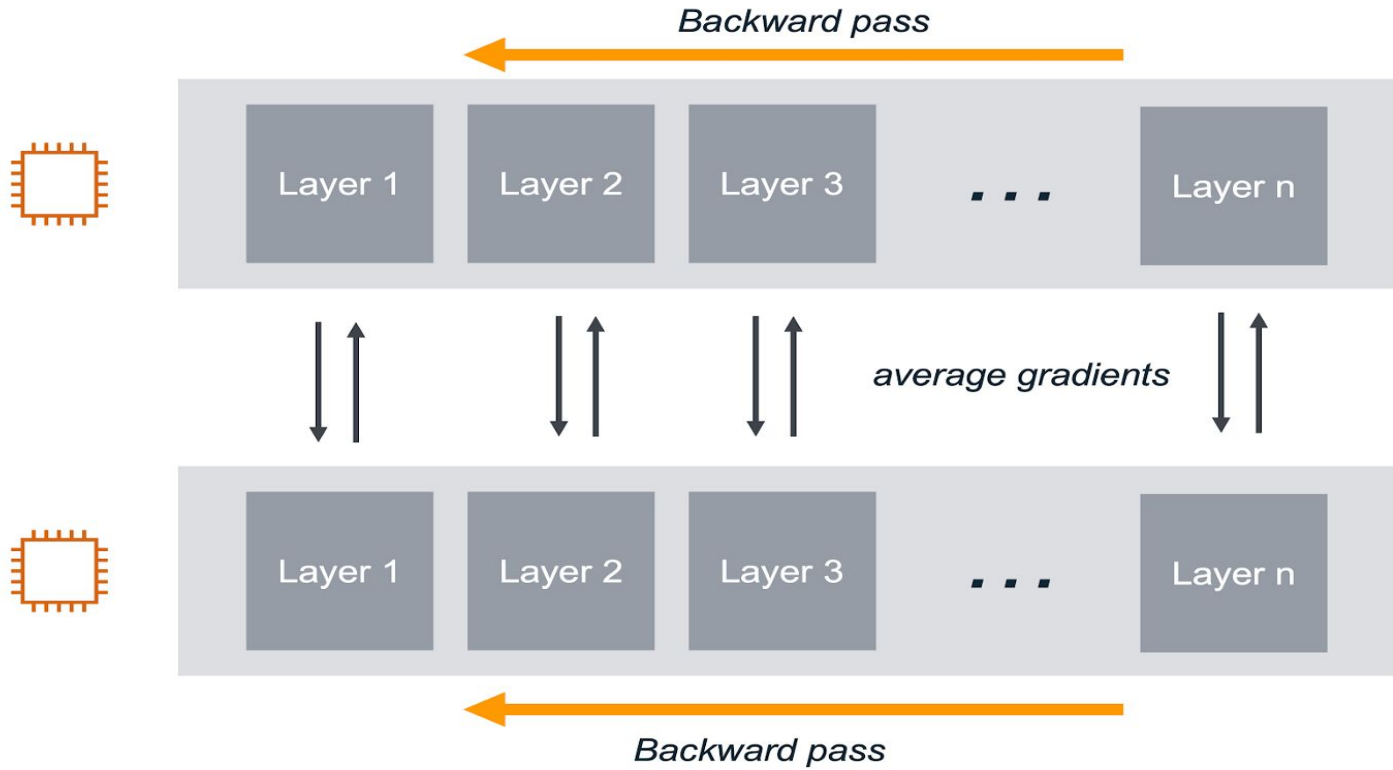
# Overall Architecture

# ML Architecture Design

# Implementation

1. Modified training script to be distributed ready using Horovod library
2. Used Horovod library to ensure the training script successfully scales to train across many GPUs in parallel, i.e. the training script is GPU count agnostic
    i.

# Dataset

Dataset: Subset of Fruits 360 from Kaggle - apples, kiwi and guava

Training images: 9600

Validation images: 1200

Testing images: 1200

# Distributed training using Amazon Sagemaker

Used Amazon Sagemaker to run the training set on scale, with the ability to scale up and down GPU's as desired

- Launched Amazon Sagemaker NB instance
- Defined the estimator with training script, location to save trained models, type of GPUs, number of GPUs per instance, TensorFlow version, MPI distribution type
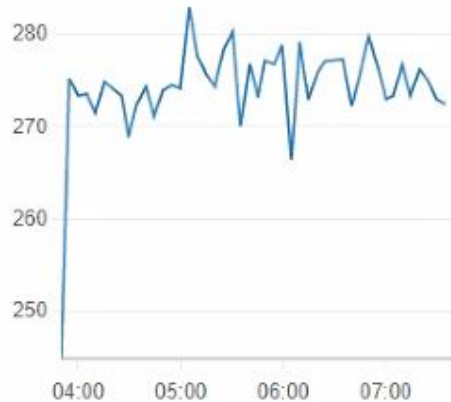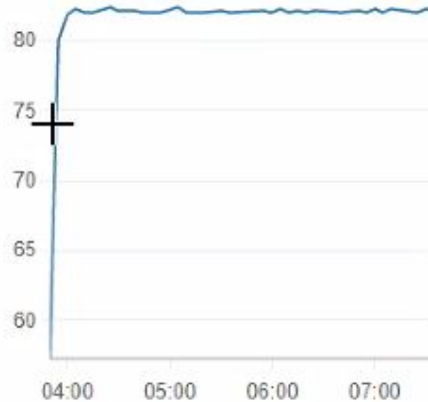- Specified paths to training, validation and test datasets in Amazon S3, passed those parameters to estimators fit function

# Monitoring
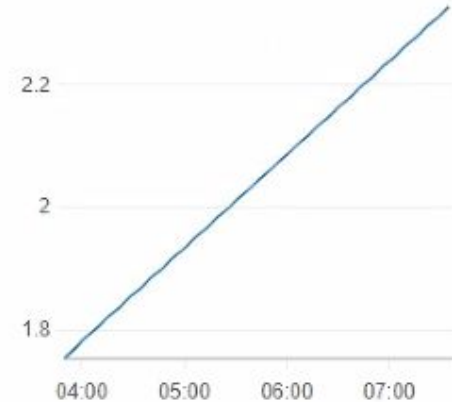
We monitored the progress through Amazon Cloud Watch

# User Interface

Web application to serve the users to use the machine learning model.
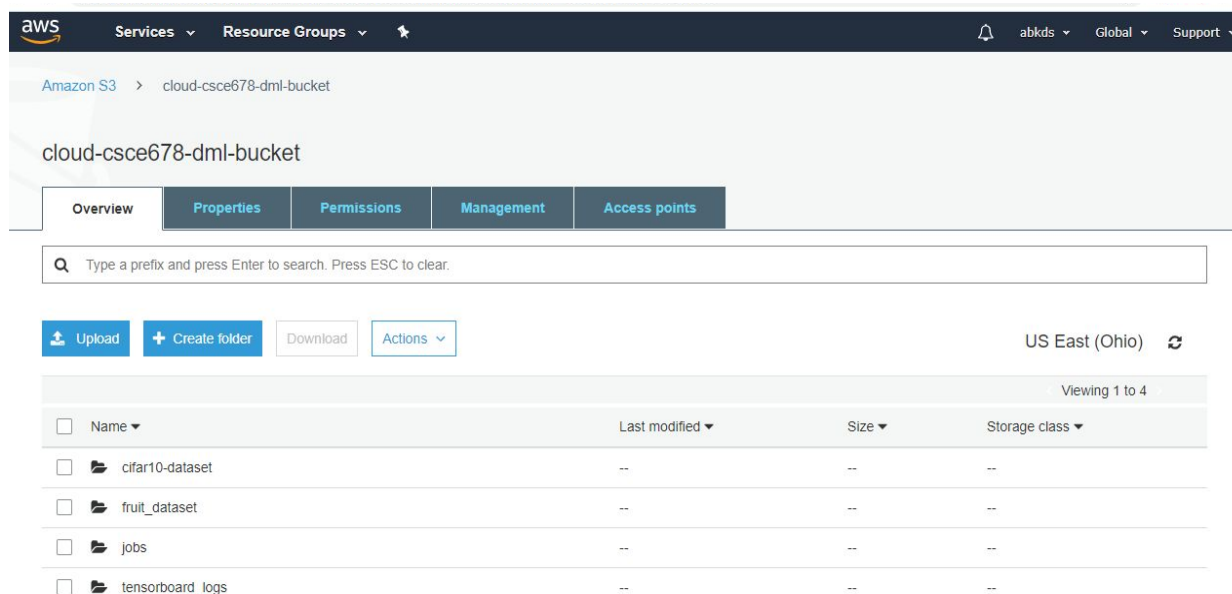
Frontend (React and Angular.js)

https://github.com/abkds/dml-frontend
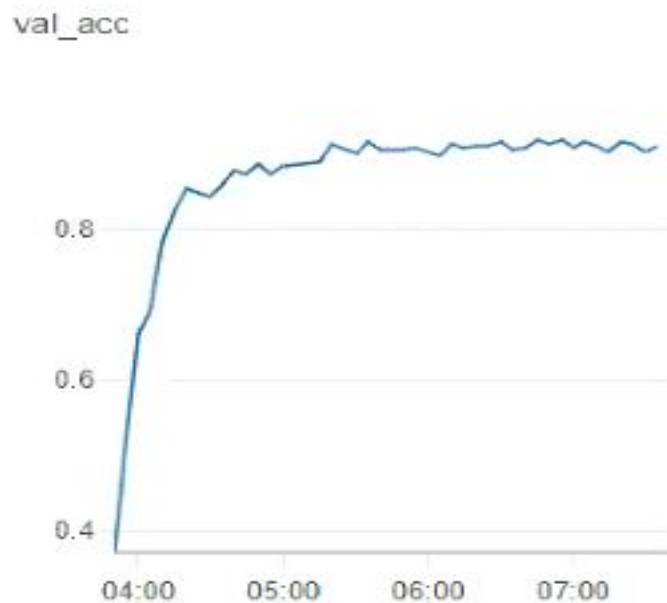
Backend: EC2 Beanstalk

# Results

Once training was complete, Sagemaker automatically uploaded training artifacts such as trained nodes, checkpoints and tensorboard logs into our S3 bucket.

# Test accuracy : 90.72



val_acc

# Thank You!