



SANTA CLARA UNIVERSITY
LEAVEY SCHOOL OF BUSINESS

STREAMLINING ML MODEL DEPLOYMENT

(Project Final Report)

For the fulfillment of the course
Software Project Management
As part of MSIS program at
Santa Clara University - Leavey School of Business

By

Ankitha Shetty

Kirtana Kirtivasan

Krishna Sai Tejaswini Kambhampati

Swetha Vijaya Raju(Team Leader)

TABLE OF CONTENTS

EXECUTIVE SUMMARY	3
LIFE CYCLE OF A ML ALGORITHM AT LEAD[2]	5
WORKFLOW OF THE ALGORITHM FROM IDEATION TO DEPLOYMENT AT LEAD	6
TEAMS & ROLES DESCRIPTION	10
CURRENT SYSTEM ARCHITECTURE	11
CURRENT OPERATIONAL COSTS[4]	13
BUSINESS REQUIREMENTS	15
FUNCTIONAL REQUIREMENTS	15
MARKET RESEARCH	16
Vertex AI[9]	16
Key Features:	16
Advantages:	17
Disadvantages:	19
VERTEX AI Demonstration	20
A. Data	20
B. Model Development	21
C. Deployment	23
MLFlow[6]	24
Advantages[7]	26
Disadvantages[7]	27
1. MLFlow Tracking API for starting and managing MLFlow runs	29
2. MLFlow Tracking UI	29
3. MLFlow Experiments	30
Competitive Analysis	33
COST COMPARISON	35
Vertex AI	35
MLFlow	37
RECOMMENDATION	38
PROCESS MODEL & PLANNING ACTIVITIES	39
Project Phases	39
Choice of Process Model	40
PROJECT PLANNING & TRACKING	41
REFERENCES	43

EXECUTIVE SUMMARY

COMPANY OVERVIEW:

Lead is a mid-size company founded in 2009 with a mission to provide a platform for communities, to not only provide real time navigation on roads, but also reports accidents, traffic jams, road works, community driven points of interests and many more activities that are the root cause of traffic congestion in cities. The company currently has 150 employees. Currently, Lead is operating only in the Bay Area, California with 2 types of users – 2000 active users who passively relay information about driving speed and traffic conditions whenever they have the application open in the background. Lead relies on large scale data solutions on a daily basis as they attempt to solve problems like - Predicting ETA, Matching Riders & Drivers and Serving the right advertisements. There are multiple engineers (Data scientists, Software engineers, reliability engineers, Data engineers and many more) working on opportunistically including ML solutions for their application.

PROBLEM SUMMARY:

Lead uses multiple ML models to explore and develop their ML algorithms for various use cases; performing semi-manual operations for training, validation and deployment; limited monitoring and validation capabilities after a model is deployed; and a hideously long deployment cycle from idea to production – This is a challenge that manifested itself causing a chaotic state at Lead. Currently, Lead does not have a service to manage the ML lifecycle.

PROJECT OBJECTIVES:

The major objective of this project is to streamline the process and minimize the timelines for deployment by reducing the manual work wherever possible in the life cycle of the algorithm, and comparatively reducing overall costs of the system using managed services. Invest in deployment processes through managed services (such as MLFlow or Vertex AI) and ensure that, by using these streamlined systems, a single Data Scientist could be able to close the product cycle from research to production.

RECOMMENDATION:*Vertex AI*

- Vertex AI provides a unified platform that can accommodate data readiness to ML model deployment streamlining the entire ML lifecycle while MLFlow will only help in Experimentation Model tracking.
- With Vertex AI, a single Data Scientist could close the product cycle from research to production.
- Vertex AI natively uses Google infrastructure and our system is majorly on Google Cloud. So the adoption of Vertex AI to our current system is compatible.

SUCCESS CRITERIA:

- The adoption of Vertex AI will bring down the operational cost of ML models in Lead by 34% in the next 5 years.
- With the adoption of Vertex AI, deployment timeline will be reduced by 22% approximately.

LIFE CYCLE OF A ML ALGORITHM AT LEAD [\[2\]](#)

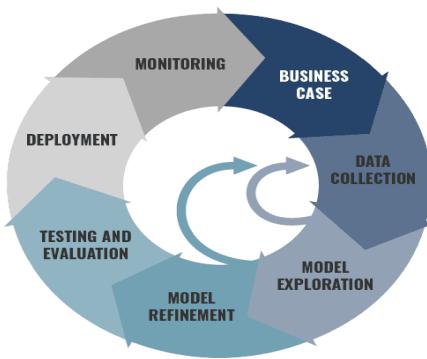


Figure 1: Life Cycle of a Machine Learning Algorithm([Image Credit](#))

All Machine Learning (ML) projects in general are highly iterative in nature; They might also need to circle back to earlier phases after completing a phase. Moreover, a project isn't complete after shipping the first version; we might have to get feedback from real-world interactions and redefine the goals for the next iteration of deployment.

Figure 1 shows the life cycle of a ML algorithm with 7 framework activities and fits them into a development model called CRoss Industry Standard Process for Data Mining (CRISP-DM), to organize projects and also use it as a reference checklist while planning and implementing them. In this process model, each phase can be iterated several times, until satisfactory results are obtained. It is also possible to jump back to the previous phase when need comes.

The seven framework activities are – Defining the task, collecting data, model exploration, model refinement, testing and evaluation, deployment and integration, monitor and maintain.

WORKFLOW OF THE ALGORITHM FROM IDEATION TO DEPLOYMENT AT LEAD

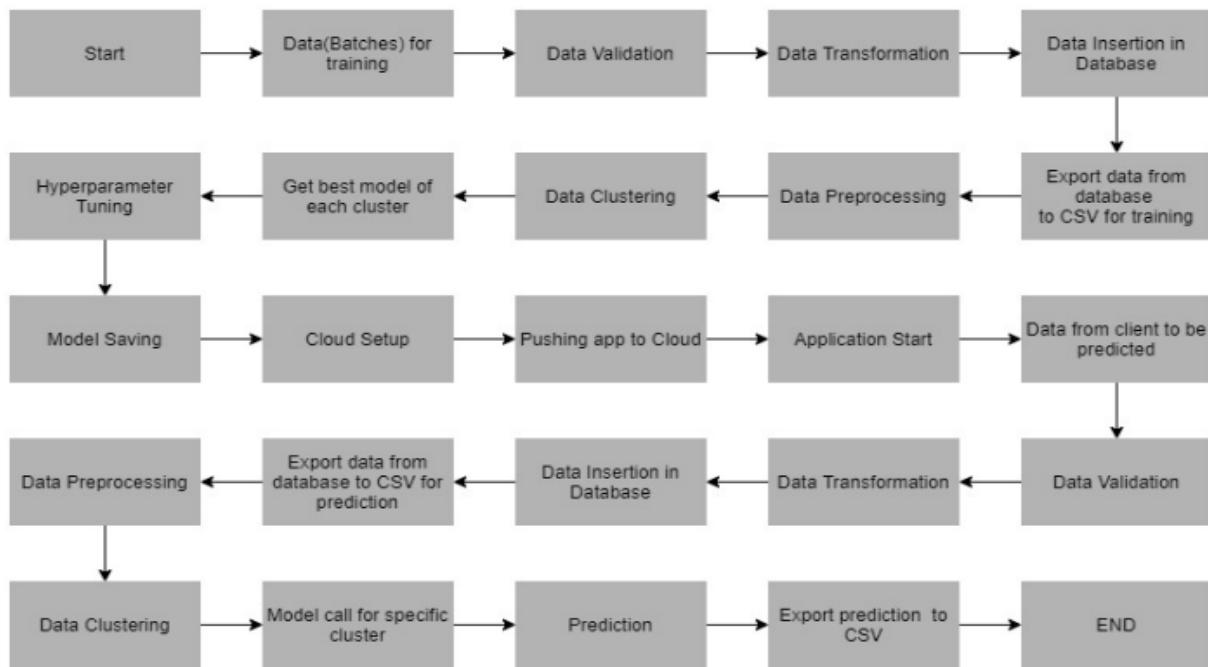


Figure 2 : Workflow of the model ([Image-credit](#))

CHECKLIST FOR MANAGING ML PROJECTS AT LEAD:

We need to understand how lead currently works to understand the problem statement better.
Lead uses the following steps to strategically plan and implement machine learning projects successfully.

Step 1 – Planning and Project Setup - There are 5 important objectives of this step (4 -6 days)

1. Defining task - **Real time ETA prediction**
2. Determining project feasibility by analyzing the below factors (10 hrs)
 - a. Cost of data acquisition for the task specified - **\$0.1/GB of data to store in BigQuery**
 - b. Cost of wrong predictions for the task - **loss of \$20 per occurrence(approx)**
 - c. Computational resources available both for training and inference
 - i. Will the model be deployed in a resource-constrained environment - **Yes**
 - ii. What are the latency requirements for the model - **Prediction latency under 10s**
3. Specifying project requirements (20 hrs)
 - a. Establishing optimization metric
4. Identifying model tradeoffs(5 hrs)
 - a. First model - Aim for a neutral first launch, focus spent on building the proper machine learning pipeline required for prediction

- b. Existing model - **Revisiting to meet updated information after discovering a promising general approach**
- 5. Setting project codebase (10 -15 hrs)
 - a. Organizing codebase to modularize data processing, model definition, model training, and experiment management - **Uses Python , Docker and Kubernetes engine**

Step 2 – Data collection and labeling – Includes 4 primary objectives (~25 hrs - 2 days)

1. Define ground truth - This task enables labeling data and ensuring ground truth is well defined as quality of the data labels has a large effect on the upper bound of a model's performance- **past eight weeks average speed of a segment on a particular time (2 hrs)**
2. Building data ingestion pipeline - **From Kafka and Google Bigquery**
3. Validating data quality (15-20 hrs)
4. Revisit step 1 and ensuring data is sufficient for the task (5 hrs)

Step 3 – Model exploration – The primary objective of model exploration is to (~5-7 days, longer for new projects)

1. Establish baselines for model performance -
 - a. First a simple model is built using initial data pipeline
 - b. This simple model is overfit to training data
 - c. Other parallel (isolated) ideas are explored at this stage
 - d. Find State of the Art model for the problem domain (if available) and reproduce results, then apply to dataset as a second baseline
 - e. Revisit Step 1 and ensure feasibility
 - f. Revisit Step a and ensure data quality is sufficient

Step 4 – Model refinement - As the team attains a general idea of successful model architectures and approaches for our use case; this step focuses on performance gains from the model with below major objectives (10 - 15 days)

1. Applying the bias variance decomposition for irreducible errors, avoidable bias, variance, and validation set overfitting (10 hrs)
2. Perform model specific optimizations (ie. hyperparameter tuning) and iteratively debug model as complexity is added (30-40 hrs)
3. Performing error analysis to uncover common failure modes - **Use clustering to uncover failure modes and improve error analysis** (25 -30 hrs)
 - a. Select all incorrect predictions. (Optionally, sorting observations by calculated loss to find the most egregious errors)
 - b. Run a clustering algorithm such as DBSCAN across selected observations
 - c. Manually explore the clusters to look for common attributes which make prediction difficult
4. Revisit Step 2 for targeted data collection and labeling of observed failure modes

Step 5 – Testing and evaluation (5 -7days)

1. Different components of a ML product to test include – (20 - 25 hrs)
 - a. Testing the full training pipeline (from raw data to trained model) which is running nightly, to ensure that changes haven't been made upstream with respect to how data from our application is stored
 - b. Running inference on the validation data and ensuring model score does not degrade with new model/weights for every code push
 - c. A quick functionality test that runs quickly (<5 minutes) ensures no broken functionality during development
 - d. Setting alerts for downtime, errors and distribution shift in data (2 hrs)

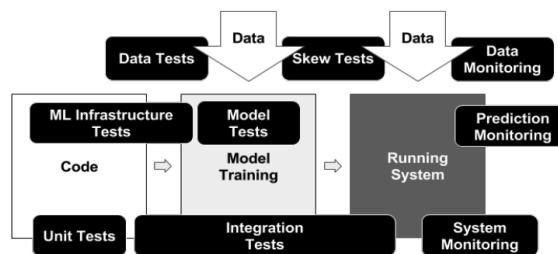


Figure 3 : ML based System Testing and Monitoring ([Image credit](#))

2. Evaluation of multiple factors for Production readiness with current System Architecture (25 hrs)
 - a. Data – Feature expectations are captured in a schema, ensuring cost of features is reasonable, Data pipeline has appropriate privacy controls, ability to add new features, all input feature code is tested
 - b. Model – Model specs are reviewed and submitted, ensure offline and online metrics correlate, all hyperparameters have been tuned, calculated impact of model staleness
 - c. Infrastructure – Ensure training reproducibility, Unit testing of model specifications, integration testing of pipeline, ensuring possibility of debugging model (if needed) before serving, model canariying² before serving, rolling back (if needed)
 - d. Monitoring – **(Currently Manual)** - Monitoring dependency changes, numerical stability, regression of computing performance and prediction quality (~20 hrs)
 - e. Revisit model evaluation metric; ensure that this metric drives desirable downstream user behavior

(superscript) 2 - a canary release is a synonym of A/B testing where a different product version is shown to a small percentage of users

Step -6 – Model deployment_[3] - Focus of the current project at Lead

1. Package the system into a Docker container and expose a REST API for inference-**Manually setting up dependencies for packaging (5 hrs)**
2. Serving the new model to small subset of users (i.e. 5%) to ensure everything goes smoothly, then roll out to all users – **By performing A/B Testing** (Measuring the delta between the new and current model's predictions will give an indication for how drastically things will change when you switch to the new model)
3. Maintain the ability to roll back model to previous versions
4. Monitoring live data (functional and operational monitoring) and model prediction distributions
- **Currently involves lot of human intervention and time (24*7 for first 15 days, once a day during peak hours for performance metrics)**

Step -7 – Ongoing model maintenance (continuous monitoring unless model is replaced)

1. Maintenance principles at Lead follows
 - a. CACE (Changing Anything Changes Everything) Principle: As ML Systems are tightly coupled, changes to any feature-space, hyper parameters, learning rate etc. can affect model performance. So lead relies heavily on model validation tests that run every time code is pushed.
2. Periodically retrain model to prevent model staleness (**happens on a weekly basis**)
3. If there is a transfer in model ownership, educate the new team

Goal: The main focus of streamlining all the steps that can be automated along the cycle by adjusting the current system architecture without compromising the principles and objectives. All such steps are identified along the process and market research is done to find alternatives.

Deployment time: With the current system architecture the Data science team working on ETA prediction needs ~ **48-50 days** to deploy the updated ML model full stream to production.

TEAMS & ROLES DESCRIPTION

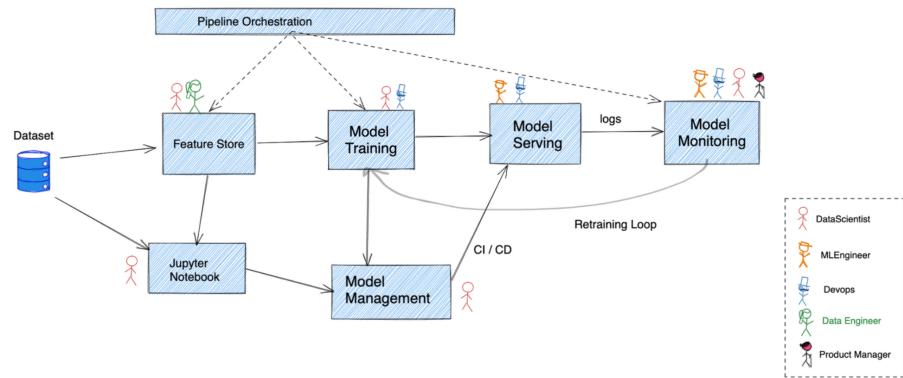


Figure 4 - Teams and roles in the overall pipeline ([Image Credit](#))

A typical Data science team at Lead comprises of the following roles –

Stakeholders	User Story
Data Engineer	ETL / Build data pipelines to load, clean and process the data
Data Scientist	Data Refinement, Data Visualization, Data Analytics
ML Engineer	Build, train and deploy the Machine Learning model
Software Engineer	Integrating the ML model with the rest of the application
Product Manager	Making product decisions and point of contact with the client

Current Data Science team working on the use case includes - 1 Data Engineer, 3 Data Scientists, 2 Software Engineers, 1 ML Engineer and a Product Manager

CURRENT SYSTEM ARCHITECTURE

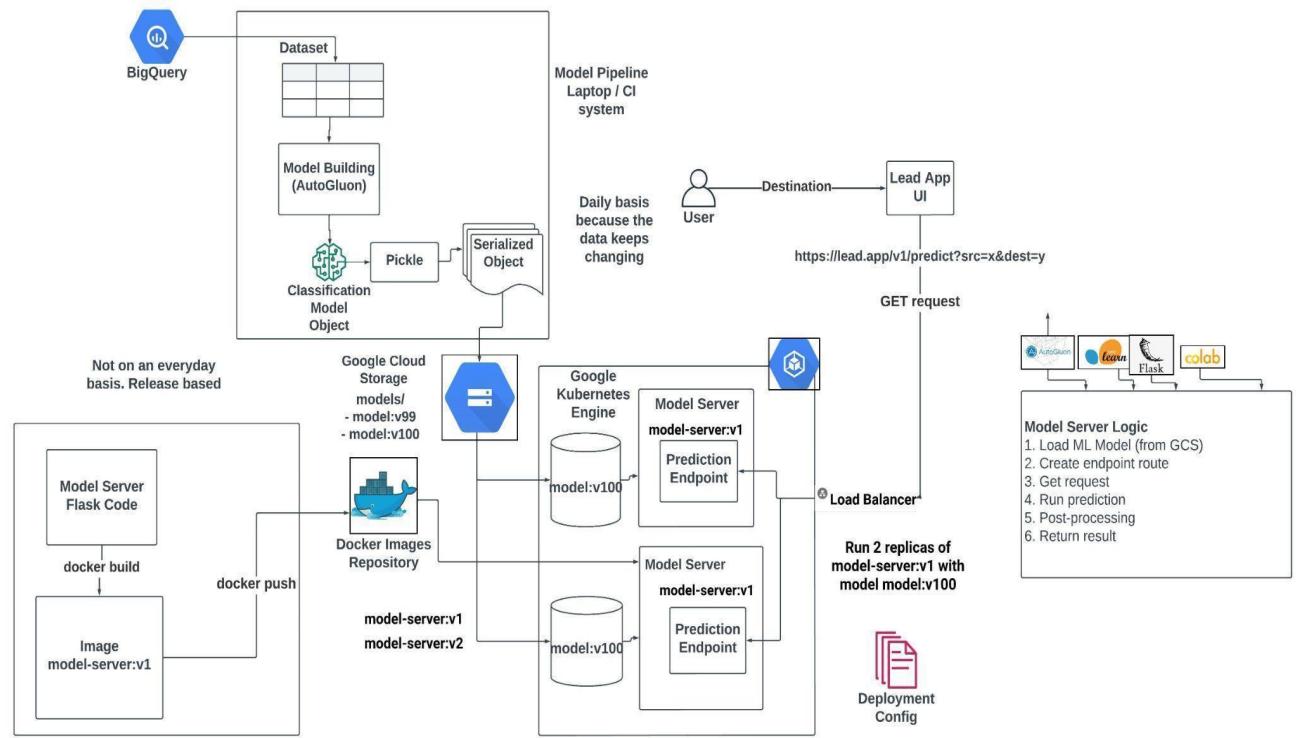


Figure 5 - Lead's current system architecture

Model Building:

Data Scientists load the datasets from BigQuery and build the model using preferred libraries (e.g. AutoGluon). The built model is serialized, versioned and stored in a GCS bucket.

Model Server:

Software Engineers parallelly build a web server with an endpoint to receive prediction results as HTTP requests. The model server is capable of loading the model stored in the GCS bucket and responds to the user results by asking the model to predict the result. The model server is built as a docker image and pushed to a docker registry (Docker Hub).

Deployment:

DevOps creates a Google Kubernetes Engine cluster with 3 nodes where Pods (containers) can be run. Once both the artifacts (model and model server) are built and pushed to their respective locations, the SREs (or Release Engineer) create a Deployment to run the model server of a specific version(e.g v1) with arguments to load a model of a specific version (e.g v100). If a new version of model server or model is to be deployed, the SREs (or Release Engineer) create a new Deployment with the new versions.

Working:

Once the model server deployment is complete, users via the Leads App can send requests to the endpoint with source and destination to get the prediction result.

Current Infrastructure

1. The instances used for serving the prediction results to users have the following configuration of at least (per instance),
 - a. CPU: 8 cores
 - b. Memory: 32GB
 - c. Disk: 2TB (enough to hold the model + overhead)
 - d. Network: 1Gbps upload/download rate
2. A load balancer to spread the load among different machines
3. The model server responds to http requests for prediction results

CURRENT OPERATIONAL COSTS [4]

Current Architecture

Training [5]

- 8 hours per model
- Per hour pricing
 - n1-standard-8, 8 CPU, 32 GB Memory = \$0.44 / hr
 - NVIDIA_TESLA_K80 GPU = \$0.617500 / hr
 - ~\$1.1 per hour = ~\$8.8 per day per model
 - 8.8*365 = \$3200 per year per model

Serving [5]

Per hour pricing

- n2-standard-8 = \$0.45 / hr
- 1 machine \$0.45 / hr = ~ \$4,000 per year (if running all the time)
- Scale up to 10 machines during a hypothetical two-hour peak that occurs everyday
 - 9 machines \$0.45 per hour = ~ \$3,000 per year (2 hrs per day)
- Total = \$7,000

Engineering / deployment

- Initial CI/CD pipeline setup cost = \$25,000

Costs for deploying 1 model to production

Model Infrastructure	Kubernetes clusters with a load balancer	\$7,000
Data Support	Independent data pipeline manager for continuous update of analytic data	\$10,000 (labor) + \$3,200
Engineering / Deployment	Continuous integration and continuous deployment (CI/CD) system to pull model from registry	\$24,500 (labor) + \$516
Total Investment / yr	\$34,500(labor) + \$11,000	

5 year TCO	\$93,500
------------	----------

Costs for deploying an additional model to production

Model Infrastructure	Cluster can be adjusted to run multiple models according to changing usage demands	+ \$3,781 / yr
Data Support	Modules can be added to pipeline management system	+ \$2,250 (labor)
Engineering / Deployment	CI/CD system can be extended to pick up additional model	+ \$2,900 (labor)
Total Investment / yr	\$5,150(labor) + \$3,781 / yr	
5 year TCO	+ \$24,055	

1. Cost for deploying 1 model to production = \$93,500 over first 5 years
2. Cost for deploying 2 models to production = \$118,000 over first 5 years
3. Cost for deploying 100 models to production = \$2.5 million over first 5 years

Goal: To comparatively reduce the overall cost of Deployment by the choice of managed service

BUSINESS REQUIREMENTS

1. The Lead application should be available 24/7
2. The prediction error percentage ³ should be less than 0.1%
3. The prediction model shouldn't return outdated information and should be retrained in case
4. Accuracy rate depends on the requirement but usually an accuracy 85% and precision above 95% is acceptable
5. The adopted managed service should be compatible with the existing system

FUNCTIONAL REQUIREMENTS

1. The deployment of a new trained model needs to be under 30 minutes
2. The model latency should be less than 5 seconds
3. The overhead latency should be less than 5 seconds (Preferably deployed to west-coast servers)
4. The prediction model should be available 24/7 with 0 downtime with at least 3 instances of the model
5. Must be able to roll back to the older models

MARKET RESEARCH

Vertex AI [\[9\]](#)

It is a managed machine learning platform offered by Google. It provides you with all of Google's cloud services in one place to deploy and maintain AI models. It helps ML workflow by providing services like: Creating a dataset and upload data, Train an ML model on your data, Upload and store your model in Vertex AI or Deploy your trained model to an endpoint for serving predictions.

Key Features:

1. Dataset: The data type can be an image, tabular, text, or video and the source can be CSV from local drive, Cloud Storage or Big query. This is very useful to import your data, connect to different datasets with GCP.
2. Feature Store: If you have created some complicated features for a dataset and are looking to build a different type of code using the same data or want to re-use the features then you can add it in the feature store and add it rather than redoing the processes again (**reduction in 20 -25 hrs of manual labor**)
3. Labeling tasks :request for a manual labeling of the images by providing the labels with correct instructions(google provides sets of guidelines for creating good instructions) in a pdf file. End user receives an email when the task is completed. You can review the progress of the data labeling task in the Google Cloud console from the Labeling tasks page. This is a paid service
4. Training: Once data for the model is uploaded you can start the training. You can train models on Vertex AI by using AutoML, or if you need the wider range of customization options available in AI Platform Training, use custom training. AutoML: Create ML models without using any code (**reduction in 5 -7 days of manual labor**)
5. Vertex AI experiments: This helps to track and analyze different run experiments. Also enables to group and compare multiple models across experiment runs

6. Metadata: Analyze ML experiments to compare the effectiveness of different sets of hyperparameters.
7. Model registry: Contains list of all the successfully trained model
8. Endpoint: components expose the functionalities of the Vertex AI.
9. Model deployment for prediction: You can deploy models on Vertex AI and get an endpoint to serve predictions on Vertex AI whether the model was trained on Vertex AI.
10. Vertex AI pipelines: Pipeline components are self-contained sets of code that perform one part of a pipeline's workflow, such as data preprocessing, data transformation, and training a model which help to streamline the development and execution of ML processes
11. Vertex AI Vizier: Used to tune hyper parameters in ML models (**3 -5 days in model refinement**)
12. Vertex AI prediction: A service optimized to run your data through hosted models
13. Use Vertex Explainable AI to get comprehensive model feature attributions and evaluation metrics.

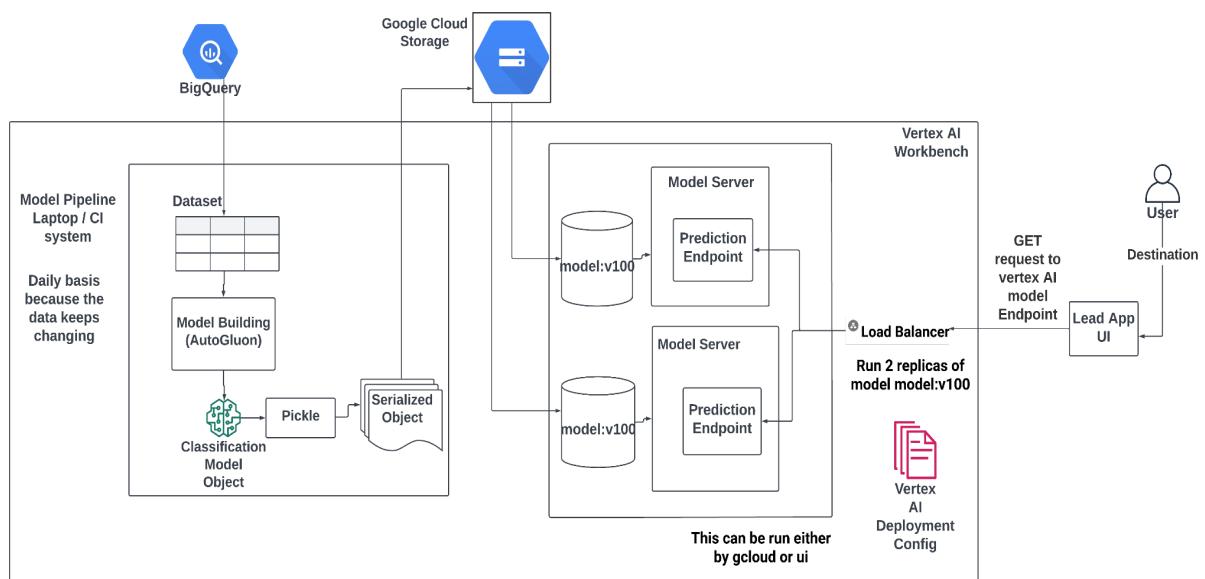


Figure 6 - Vertex AI integration with the existing architecture

Advantages:

1. Easy to access, takes as low as 30 seconds to login and access Vertex AI
2. New customers get \$300 in free credits to spend on Vertex AI.
3. Vertex AI provides the user with a dashboard which has consolidated details of the project
4. Data can be uploaded to the cloud quickly. Eg. 1 million records were uploaded in less than a minute. The uploaded data is stored on google cloud bucket.
5. Provides AutoML models that will predict/classify the data provided without any major cleanup. If the requirements are simple classification and team size is small and they do not specialize in ML then AutoML is the best choice. AutoML can also be used to create an initial proof of concept (prototype) or baseline model.
6. Very intuitive environment, training the data was very simple, and the user did not have to write any code. Users can simply select the column to predict and put the data.
7. Vertex AI monitors the training and notifies the users via email when there is a failure or success of model training. The message contains the error message for further analysis.
8. After training the data, Vertex AI provides the performance details of the trained data using confusion matrix and the losses.
9. The code can be integrated with Github.
10. User customized code (.ipynb) can be easily added to Vertex AI workbench by simple drag and drop. In custom training, you can select from among many different machine types to power your training jobs
11. Vertex AI provides ML monitoring capable of tracking metrics, boosting the observability of your project
12. Creating an endpoint is very simple. Steps are very clear and an endpoint is created in a couple of minutes. This can be used for batch and online prediction options.
13. Google cloud provides you with a [cost management tool](#). This will help the administrator (project manager) to assign a budget for the project. It also provides a dashboard to summarize the spending. There are features to alert the users through emails.

14. In Depth training [videos](#) and codes on how to set up and run ML models on vertex this would be helpful since we would have to spend on external training. But the team's time will have to be allocated for the said training.

Disadvantages:

1. Vertex AI can become very costly if AutoML is used for model training. For tabular data[\[5\]](#),

Operation	Price per node hour (classification/regression)
Training	\$ 21.252
Deployment & online prediction	Same as custom model (\$0.39)

2. Vertex AI natively uses infrastructure from Google. The team members should be trained to use Vertex AI from data readiness to deployment.
3. The data scientists who will be using Vertex AI should have some familiarity about the cloud technologies like containerising, cloud buckets, etc.

VERTEX AI Demonstration

The screenshot shows the Vertex AI Dashboard with several sections highlighted:

- Recent datasets:** Lists datasets like Traffic_dataset, Flowers, Source_Big_query, and Traffic_unclean.
- Recent models:** Lists models like Flowers (Version 1), traffic_50M_classify (Version 1), Custom_mpg_model (Version 1), and Predict_Traffic_type (Version 1).
- Recent batch predictions:** Shows a prediction for PredictType.
- Recent notebook instances:** Lists instances like experiment-20221118-121553, traffic-predict, and tensorflow-2-6-20221116-104849.
- Recent errors:** Lists errors from Custom_ML_model-custom-job and Custom_ML-custom-job.

A callout box notes: "The dashboard gives visual display of all tasks that we have recently worked on. The primary intention is to provide information at-a-glance. Here we can see the recent datasets, models, predictions, notebook instances and recent errors."

[1] Vertex AI DashBoard

A. Data

The screenshot shows the Vertex AI Datasets page for the Source_Big_query dataset:

- Properties:** Created: Nov 17, 2022 11:29 AM, Dataset format: BigQuery, Dataset location(s): bq://dependable-gl-dataset.Traffic_CAE, Encryption type: Google-managed key.
- Summary:** Total columns: 15, Total rows: 300,000.
- Statistics:** Column types and distribution: STRING (9 (60%)), INTEGER (1 (6.67%)), TIMESTAMP (2 (13.33%)), FLOAT (3 (20%)).
- General statistics:** Generated by Nov 17, 2022 11:58 AM, with a link to GENERATE STATISTICS.

Callouts explain:

- Source: If you want to change the source of the data, you can do so by simply clicking on the source and choose the new data file. Then the new data will be stored in cloud storage again. Changing the source does not mean that the previous file in the cloud storage is deleted.
- Data type: The data type can be an image, tabular, text, or video and the source can be CSV from local drive, Cloud Storage or Big query. This is very useful to import your data, connect to different dataset within GCP. If data is loaded from the local computer it is then stored in the google cloud bucket (click on the path).
- Statistics: Once the data is uploaded successfully you can generate the statistics for the data. This gives information about the number of rows , columns, distinct values in each column, missing value % in each column.

[2] Vertex AI Datasets

The screenshot shows the Vertex AI Labeling Task page for the Flowers dataset:

- Dataset summary:** IMPORT, BROWSE, ANALYZE, Filter items: Select all.
- Training data:** All (39), Labeled (39), Unlabeled (0), Training (30), Validation (5), Test (4).
- Labels:** Flower (29), Not_Flower (10). A green box highlights "ADD NEW LABEL" and a green box labeled "Labels".
- Image preview:** A grid of flower images with labels: Not_Flower, Flower, Not_Flower, Flower, Not_Flower, Flower, Flower, Flower, Flower, Flower.
- Jobs and models:** Training jobs and models for Flowers (Version 1) with Model type: Image classification, Resume Training, TRAIN NEW MODEL.
- Labeling tasks:** A green box highlights CREATE LABELING TASK.
- Instructions:** In the image data, you can upload up to 500 images per upload. You can choose the file and select data split. The data split has the option of Training, test, or validation. So, the data can be split into groups during the upload. The google cloud storage path where the images are uploaded will be stored in this path. Once the image is uploaded, you can manually label the image or you can create a labeling task. A labeling task is to request for a manual labeling of the images by providing the labels with correct instructions (google provides sets of guidelines for creating good instructions) in a pdf file. Users receives an email when the task is completed. You can review the progress of the data labeling task in the Google Cloud console from the Labeling tasks page.

[3] Vertex AI Labeling Task

B. Model Development

Vertex AI training tab

Training pipelines are the primary model training workflow in Vertex AI. You can use training pipelines to create an AutoML-trained model or a custom-trained model. For custom-trained models, training pipelines orchestrate custom training jobs and hyperparameter tuning with additional steps like adding a dataset or uploading the model to Vertex AI for prediction serving. [Learn More](#)

Once data for model is uploaded you can start the training. You can train models on Vertex AI by using AutoML, or if you need the wider range of customization options available in AI Platform Training, use custom training. To create a training job click on create option. Training page gives the details of training like status. Once the job completes user is notified with an email.

Name	ID	Status	Job type	Model type
Flowers-20221118-130156	5913045715952599040	Finished	Training pipeline	Image classification (Single-label)
Flowers	2498771840638386176	Finished	Training pipeline	Image classification (Single-label)
traffic_50M_classify-20221117-114717	1228475270743195648	Finished	Training pipeline	Tabular classification
Custom_mpg_model	8205993652795736064	Finished	Training pipeline	Custom
Custom_ML_model	831745087162548224	Failed	Training pipeline	Custom

Image:[4]

Steps to create an AutoML model in Vertex AI[5]

Step 1: Train new model

1 Training method
2 Model details
3 Training options
4 Compute and pricing

Dataset: Traffic_dataset
Objective: Classification

Model training method:
 AutoML
 Train high-quality models with minimal effort and machine learning expertise. Just specify how long you want to train. [Learn more](#)
 Custom training (advanced)
 Run your TensorFlow, scikit-learn, and XGBoost training applications in the cloud. Train with one of Google Cloud's pre-built containers or use your own. [Learn more](#)

Step 2: Train new model

Name: Traffic_dataset
Description
Target column: Type

Data split:
 Random assignment
 80% of your data is randomly assigned for training, 10% for validation, and 10% for testing.
 Training: 80%, Validation: 10%, Test: 10%
 Default: 0%
 Manual
 You assign each data row for training, validation, and testing. [Learn more](#)
 Chronological assignment
 The earliest 80% of your data is assigned to training, the next 10% for validation and the latest 10% for testing. This option requires a Time column in your dataset. [Learn more](#)

Step 3: Train new model

1 Before continuing, use the Transformation column to review and specify the data types in your dataset. If unspe

Column name	Transformation	Missing % (count)	Distinct values	Ca
Address	Categorical	95	-	-
Age	Numeric	752	-	-
City	Text	39	-	-
County	Text	521802	-	-
Description	Text	4251	-	-
Distance	Numeric	261101	-	-
EndTime	Text	138399	-	-
LocationLat	Numeric	134311	-	-
LocationLong	Numeric	5	-	-
Severity	Categorical	3	-	-
Side	Text	285674	-	-
StartTime	Text	15520	-	-
Street	Text	2	-	-
TimeZone	Categorical	7	-	-
Type	Categorical	-	-	-

Total 14 feature columns are included in the training

Weight column
Select a column to specify how to weight each row of the training data. By default, each row of your training data is weighted equally. [Learn more](#)

Optimization objective:
 AUC ROC
 Discretize between classes
 Log loss
 Keep prediction probabilities as accurate as possible
 AUC PR
 Maximize precision-recall for the less common class
 Precision at recall
 Maximize precision for the less common class
 Recall at precision
 Maximize recall for the less common class

Step 4: Train new model

1 Enter the maximum number of node hours you want to spend training your model.
You can train for as little as 1 node hour. You may also be eligible to train with free node hours. [Pricing guide](#)

Budget: Maximum node hours: Estimated completion: Enable early stopping

Click on create and chose the dataset to work on and chose AutoML. In the model details chose if you are training a new model or a new version and then chose the target column. You can click on advanced and chose the data splits type. Your dataset will be split into training, validation, and testing sets. The default split Vertex AI applies depends on the type of model that you are training.

Next in training Options we can see the data schema, we can manually choose the datatype or else AutoML chooses it for us. You can then chose the classification metrics. Next chose the budget, the budget refers to the number of node hours (represents the time a virtual machine spends running your training) spent on training the model. For each hour vertex ai charges 21 dollars.

Image:[5]

Custom Training:

Workbench[6]

The screenshot shows the Vertex AI Workbench interface. On the left, there's a sidebar with various tools like Dashboard, Workbench (which is selected and highlighted with a red box), Pipelines, Feature Store, Datasets, Labeling tasks, Training, Experiments, Metadata, Model Registry, Endpoints, Batch predictions, and Matching Engine. Below these are Marketplace and other sections.

The main area is titled "Create a user-managed notebook". It has fields for "Notebook name" (set to "python-20221121-17465a"), "Region" (set to "us-west1 (Oregon)"), and "Zone" (set to "us-west1-b"). A note says: "Requests to your notebook from the DataLab/Jupyter interface may be routed through a different region than selected above depending on service availability."

Under "Environment", "Operating System" is set to "Debian 10" and "Environment" is set to "Python 3 (with Intel® MKL)". A note states: "All environments have the latest NVIDIA GPU libraries, Intel libraries and drivers. Notebooks use JupyterLab 3 by default (you can specify a previous version instead)."

There's a "Custom metadata" section with a "+ ADD ITEM" button and a "Metadata" section with a "+ ADD ITEM" button.

Under "Machine configuration", "Machine type" is set to "n1-standard-4 (4 vCPUs, 15 GB RAM)" (highlighted with a red box). A note says: "GPU type None". Another note at the bottom says: "There are no GPUs available for the zone, environment, and machine type selected above. Learn more".

On the right, there's a "DETAILS" section with a note: "To create your custom code, go to workbench tab. Vertex AI Workbench is a Jupyter notebook-based development environment for the entire data science workflow. Both notebook options are prepackaged with JupyterLab and have a preinstalled suite of deep learning packages, including support for the TensorFlow and PyTorch frameworks. You can specify the name, location, use CPU-only or GPU-enabled instances. Your Vertex AI Workbench notebook instances are protected by Google Cloud authentication and authorization."

Image:[6]

Open Jupyter lab[7]

The screenshot shows the Google Cloud Platform interface. On the left, there's a sidebar with "Vertex AI" (selected and highlighted with a yellow box), "Dashboard", "Workbench" (highlighted with a red box), "Pipelines", "Feature Store", "Datasets", "Labeling tasks", "Training", "Experiments", and "Metadata".

The main area shows "MANAGED NOTEBOOKS" and "USER-MANAGED NOTEBOOKS". Under "USER-MANAGED NOTEBOOKS", there are two entries: "experiment-20221118-121553" and "my-notebook" (highlighted with a red box). Both entries have "OPEN JUPYTERLAB" next to them.

To the right, there's a "Notebook details" panel for "my-notebook". It shows "BASIC INFO" with the following details:

- Region: us-central1 (Iowa)
- Zone: us-central1-a
- Environment: Python 3 (with Intel® MKL)
- Environment version: M100
- Machine type: n1-standard-4 (4 vCPUs, 15 GB RAM)
- GPU: None
- Boot disk: 100 GB disk
- Data disk: 100 GB disk
- Created: Nov 15, 2022, 3:58:25 PM
- Last modified: Nov 15, 2022, 3:59:19 PM
- Backup: default
- Subnetwork: 438574652846-compute@developer.gserviceaccount.com
- Service account: Service account
- Permission mode: Service account

A "OPEN JUPYTERLAB" button is highlighted with a red box.

Below this, there's a Jupyter Lab interface with a terminal tab (Terminal 1) and a code tab (Code). The terminal tab shows the command "SPM.ipynb" (highlighted with a red box). The code tab shows some Python code related to data processing and storage.

On the far right, there's a vertical toolbar with various Git-related options: Initialize a Repository, Clone a Repository, Merge Branch..., Push to Remote, Push to Remote (advanced), Pull from Remote, Pull from Remote (Force), Reset to Remote, Manage Remote Repositories, Open Git Repository in Terminal, Simple staging, and Double click opens diff.

Image:[7]

Vertex AI Experiments tab[8]

The screenshot shows the Vertex AI interface with the 'Experiments' tab selected. On the left sidebar, under 'MODEL DEVELOPMENT', 'Experiments' is highlighted. The main area displays a table titled 'Runs' with one entry: 'run-1'. A callout box highlights the text: 'Vertex AI Experiments This helps to track and analyze different run experiments. This is important when comparing models or performing hyper parameter tuning'.

Image:[8]

Model Registry: Vertex Explainable AI on evaluation metrics[9]

The screenshot shows the Vertex AI interface with the 'Model Registry' tab selected. Under 'EVALUATE', a model named 'traffic_50M_classify' is selected, and its 'Version 1' is shown. The 'Evaluation' section displays various metrics: PR AUC (1), ROC AUC (1), Log loss (0.018), F1 score (0.9929882), Precision (99.3%), and Recall (99.3%). Below these metrics are three plots: 'Precision-recall curve', 'ROC curve', and 'Precision-recall by threshold'. A callout box highlights the 'Confusion matrix' section, which provides a heatmap showing the percentage of correct classifications and misclassifications for each label (Congestion, Accident, Flow-Incident, Lane-Blocked, Broken-Vehicle) against each predicted label.

Image:[9]

C. Deployment

Endpoint [10]

The screenshot shows the Vertex AI interface with the 'Endpoints' tab selected. The table lists endpoints, including 'v1_custom_mpg' and 'Predict_trafficType_endpoint'. An orange arrow points from the text 'Endpoint is component that expose the functionalities of the Vertex AI' to the 'Predict_trafficType_endpoint' row. The table columns include Name, ID, Status, Models, Region, and Monitoring.

Image[10]

MLFlow_[6]

- MLflow is an open source platform for managing the end-to-end machine learning lifecycle
- Some benefits of MLFlow are Experiment Tracking, Model Management, and Model Deployment
- Four primary components of MLFlow are: Tracking, Projects, Models and Model Registry

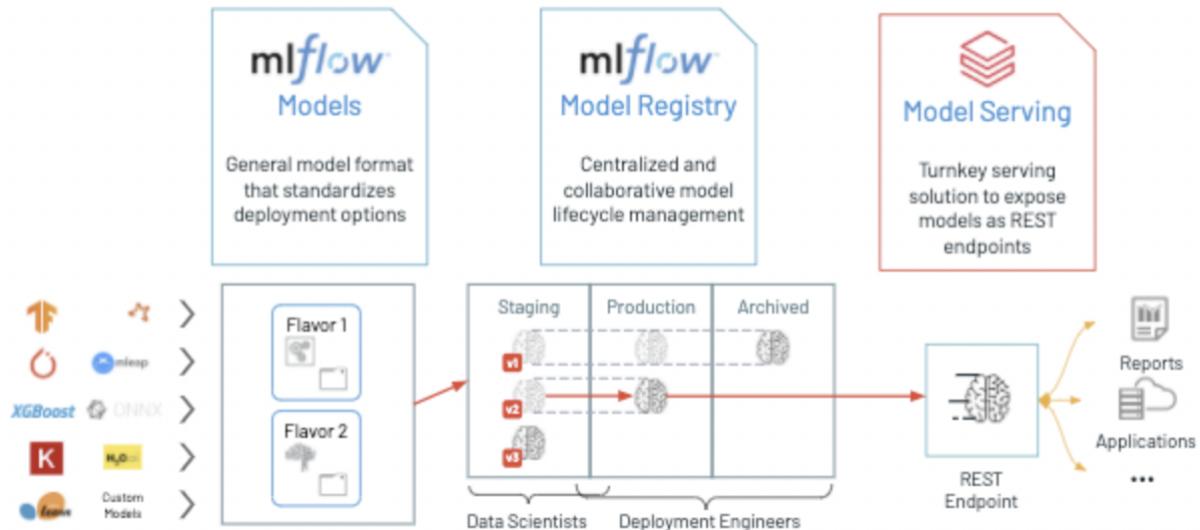


Figure 7 - MLFlow Serving ([Image Credit](#))

MLFlow Components

Tracking (25-30 hrs of manual work)

- MLflow Tracking is organized around the concept of runs, which are executions of ML science code
- MLflow tracking provides an API and UI for logging parameters, code versions, metrics, and artifacts when running code and for later visualizing the results
- MLflow allows you to group runs under experiments, which can be useful for comparing runs intended to tackle a particular task
- You can record runs using MLflow Python, R, Java, and REST APIs from anywhere you run your code
- The MLFlow runs can be recorded either to local files, to a SQLAlchemy compatible database, or remotely to tracking server
- An MLFlow tracking server has two components: a backend store and an artifact store

- The backend store is where MLflow Tracking Server stores experiment and run metadata as well as parameters, metrics, and tags for runs
- The artifact store is a location suitable for large data (such as the GCS bucket) and is where you log their artifact output (for example, models)

Projects

- An MLflow Project is a format for packaging data science code in a reusable and reproducible way
- Each project is a directory of files, or a Git repository, containing your ML code. For example, projects can contain a `conda.yaml` file for specifying a Python Conda environment
- Each Project specifies properties such as: a human-readable *Name* of the project; *Entry Points* which include commands that can be run within the project and information about their parameters (calling `.py` file in project as entry point); *Environment* is the software environment that should be used to execute project entry points (includes all dependencies required by the project code)
- Software Environments supported by MLFlow Projects include Conda environments, Virtualenv environments, and Docker containers

Models (**2-4 days of manual work**)

- An MLflow Model is a standard format for packaging machine learning models that can be used in a variety of downstream tools (for e.g., real-time serving through a REST API or batch inference on Apache Spark)
- Offers a convention for packaging machine learning models in multiple flavors, and a variety of tools to help you deploy them
- MLflow provides tools to deploy many common model types to diverse platforms: for example, any model supporting the “Python function” flavor can be deployed to a Docker-based REST server, to cloud platforms such as Azure ML and AWS SageMaker, and as a user-defined function in Apache Spark for batch and streaming inference.

Registry (~3 days of manual work)

- Offers a centralized model store, set of APIs, and UI, to collaboratively manage the full lifecycle of an MLflow Model
- It provides model lineage (which MLflow experiment and run produced the model), model versioning, stage transitions (for example from staging to production or archiving), and annotations

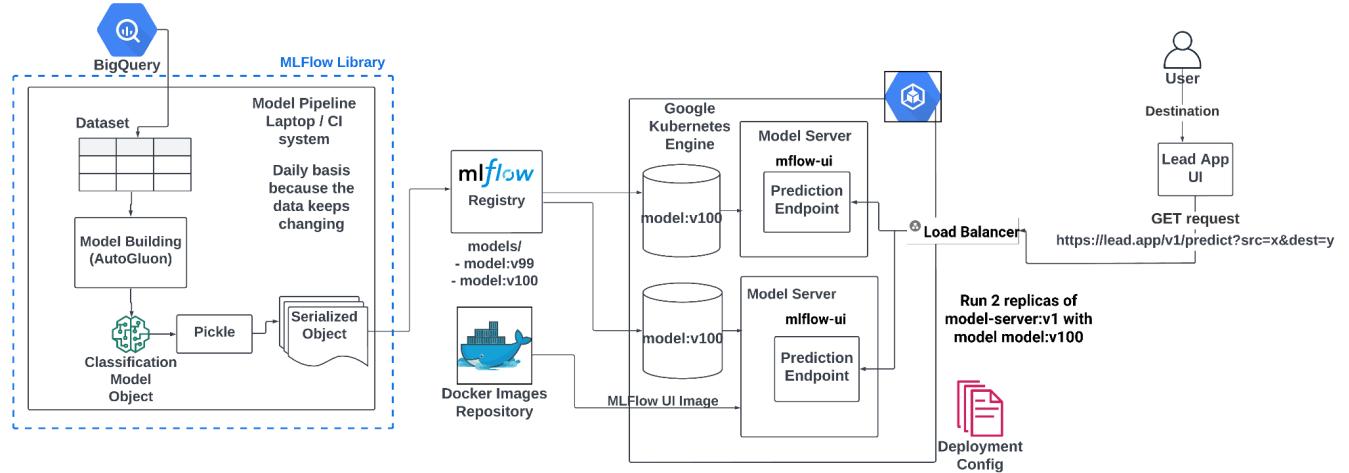


Figure 8 - MLFlow integration with the existing architecture

Advantages [\[7\]](#)

1. MLFlow is an Open Source MLOps tool and so all the features offered are free.
2. Offers model tracking, management, packaging, and centralized lifecycle stage transitions.
3. Aids in creating a collaborative environment for model development.
4. It is easy to adapt ML experiments to MLFlow as the tracking is done via a simple import in your code; It is a Python package that enables experiment tracking on current ML algorithms.
5. Experiment tracking, which is the core of MLFlow, favors the ability to develop locally and track runs in remote locations via a logging process. This is suitable for exploratory data analysis, model evaluation, and tuning hyperparameters in the ML Lifecycle.
6. MLFlow achieves model deployment by utilizing the Model Registry. Through this, MLFlow provides a central location to share ML models as well as a space for collaboration on how to move them forward for deployment.

7. Model Registry provides model versioning, model lineage, annotation, and stage transitions. It has a set of APIs and a UI to manage the complete lifecycle of the MLflow model more collaboratively.
8. MLflow makes it easy to promote models on different cloud environments like Amazon Sagemaker, Microsoft Azure, and Google Cloud Platform.
9. Meets the needs of data scientists looking to organize themselves better around experiments and machine learning models. For such teams, ease of use and setup is often the main driver.
10. Saves time for data scientists as MLflow provides features to visualize and compare experiments which contain various versions of the ML model performing a particular task (say, regression or classification).
11. The parameters for the ML model can be supplied to the pipeline as arguments which saves time as compared to re-running the entire Jupyter notebook. It also lets data scientists compare the performance of the model with changing parameters by logging all runs in the MLflow Tracking server.
12. MLflow provides out of the box APIs for Python, R, and Java. For any other language, it also provides REST APIs which allow the user to log runs, metrics, and parameters seamlessly.
13. MLflow has provided [Documentation](#) which gives in-depth information on all their components, how to install, and various samples to implement MLflow.

Disadvantages [\[7\]](#)

1. MLflow requires the data scientist to perform data cleaning and data pre-processing steps as well as model development from scratch rather than MLflow providing an application to perform these steps; This requires the user to have prior ML model development knowledge/experience.
2. Does not automatically find a model that gives the best results (AutoML not supported).
3. Since AutoML is not supported, model development might take several days to weeks just to target a specific application for Lead.

4. MLFlow requires the current architecture for the deployment to be complete. MLFlow Recipes (also known as MLFlow Pipelines), a framework that enables data scientists to quickly develop high-quality models and deploy them to production, is still in the experimental phase.
5. The cost does not reduce much even though this is an open source tool as it requires the presence of current architecture for deployment and it still does not provide automation for ML model development unlike Google's Vertex AI.

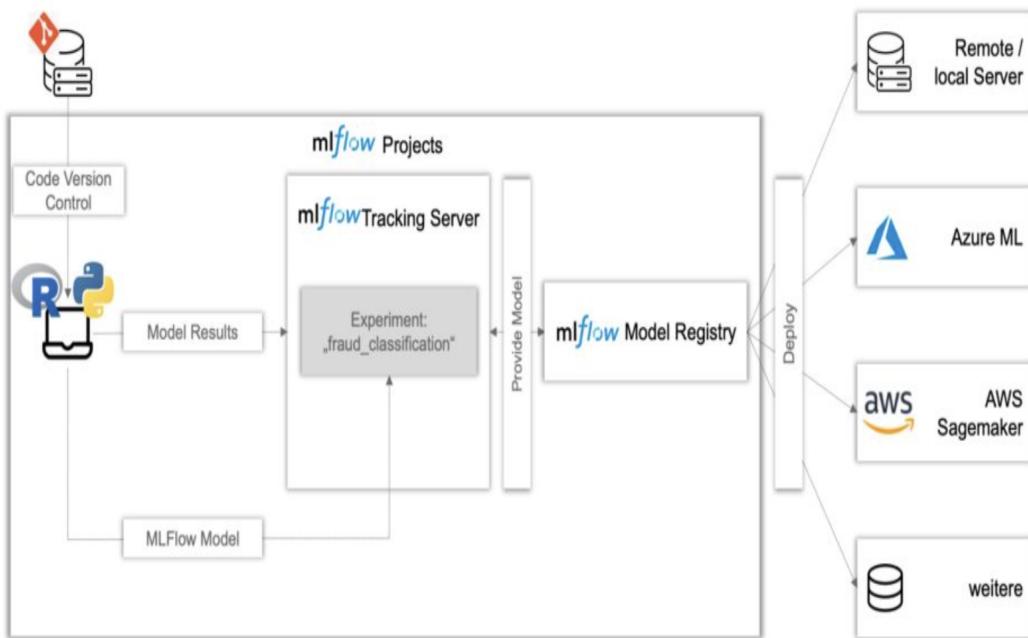


Figure 10 - Workflow Management with MLFlow ([Image Credit](#))

MLFLOW DEMONSTRATION

1. MLFlow Tracking API for starting and managing MLFlow runs

```
import mlflow
import mlflow.sklearn
import os

os.environ["MLFLOW_TRACKING_URI"] = "http://127.0.0.1:5000/"

eval_data = X_test
eval_data["label"] = y_test

mlflow.set_experiment(experiment_id = '0')

with mlflow.start_run() as run:
    # Log the baseline model to MLflow
    mlflow.sklearn.log_model(model, "model")
    model_uri = mlflow.get_artifact_uri("model")

    # Evaluate the logged model
    result = mlflow.evaluate(
        model_uri,
        eval_data,
        targets="label",
        model_type="classifier",
        evaluators=["default"],
    )
```

2. MLFlow Tracking UI

The screenshot shows the MLflow 2.0.1 tracking interface. At the top, there's a navigation bar with 'mlflow' logo, 'Experiments', 'Models', 'GitHub', and 'Docs'. Below the navigation is a search bar and a dropdown menu for 'Default' experiment. The main area displays a table of 'Matching runs' with columns: Run Name, Created, Duration, Source, Models, and Metrics. The Metrics column lists accuracy_sc, f1_score, log_loss, mae, precision_reca, and prec. There are 6 rows in the table, each corresponding to a different run name and its details.

Run Name	Created	Duration	Source	Models	Metrics
powerful-gull-150	3 days ago	3.2s	Logistic...	sklearn	accuracy_sc: 0.886, f1_score: 0.912, log_loss: 0.335, mae: -, precision_reca: 0.932, prec: 0.89
trusting-colt-811	7 days ago	20.0min	MLFlow...	shap, 1 more	f1_score: 1, log_loss: 3.374e-4, mae: -, precision_reca: 1, prec: 1
judicious-elk-333	7 days ago	3.1s	MLFlow...	sklearn	accuracy_sc: -, f1_score: -, log_loss: -, mae: 0.517, precision_reca: -, prec: -
worried-wren-72	7 days ago	3.2s	MLFlow...	sklearn	accuracy_sc: -, f1_score: -, log_loss: -, mae: 0.563, precision_reca: -, prec: -
agreeable-eel-729	7 days ago	2.9s	MLFlow...	sklearn	accuracy_sc: -, f1_score: -, log_loss: -, mae: -, precision_reca: -, prec: -
caring-croc-360	7 days ago	3.1s	MLFlow...	sklearn	accuracy_sc: -, f1_score: -, log_loss: -, mae: 0.627, precision_reca: -, prec: -

3. MLFlow Experiments

The screenshot shows the MLFlow UI interface. At the top, there's a navigation bar with 'Experiments' and a dropdown menu. Below it is a search bar with 'Default' selected. The main area displays a table of 'Matching runs' with columns: Run Name, Created, Duration, Source, Models, and Metrics. The Metrics column includes accuracy_score, f1_score, and log_loss. A circled 'Default' entry is highlighted.

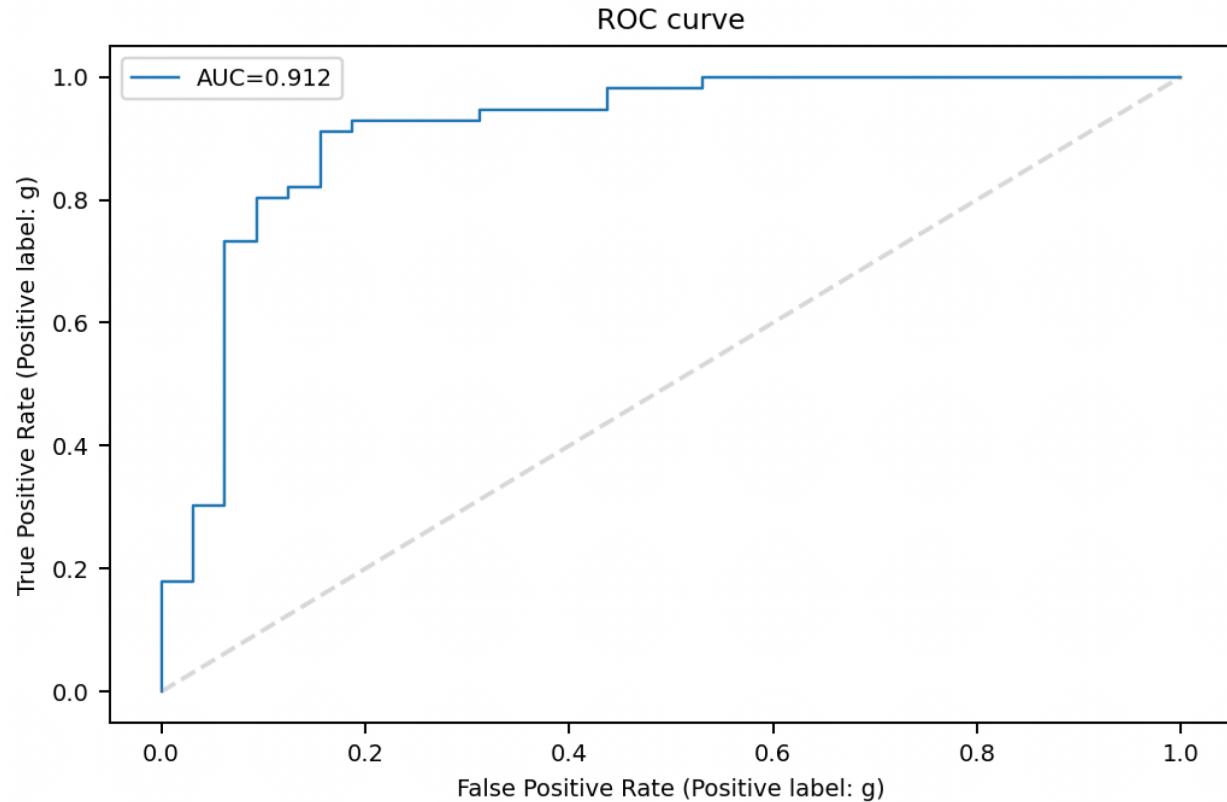
	Run Name	Created	Duration	Source	Models	Metrics
powerful-gull-150	3 days ago	3.2s	MLFlow...	sklearn	accuracy_score: 0.886 f1_score: 0.912 log_loss: 0.335	
trusting-colt-811	7 days ago	20.0min	MLFlow...	shap, 1 more	f1_score: 1 log_loss: 3.374e-4	
judicious-elk-333	7 days ago	3.1s	MLFlow...	sklearn	-	
worried-wren-72	7 days ago	3.2s	MLFlow...	sklearn	-	
agreeable-eel-729	7 days ago	2.9s	MLFlow...	sklearn	-	
caring-croc-360	7 days ago	3.1s	MLFlow...	sklearn	-	

4. MLFlow Logged Metrics

Metrics (13)

Name	Value
accuracy_score ↗	0.886
example_count ↗	88
f1_score ↗	0.912
false_negatives ↗	4
false_positives ↗	6
log_loss ↗	0.335
precision_recall_auc ↗	0.932
precision_score ↗	0.897
recall_score ↗	0.929
roc_auc ↗	0.912
score ↗	0.886
true_negatives ↗	26
true_positives ↗	52

5. MLFlow Visualization



6. Filter Runs according to Metrics/Parameters/Attributes

The screenshot shows the MLflow 2.0.1 interface. The top navigation bar includes "Experiments" and "Models". The main area displays the "Default" experiment with an ID of 0 and artifact location. It shows a list of 6 matching runs, each with a checkbox, run name, creation date, recall score, and RMSE. A search bar filters for runs with RMSE less than 1 and using a tree model. On the right, a sidebar provides filtering options for "Attributes" (User, Source, Version, Models) and "Metrics" (accuracy_score, f1_score, log_loss, mae, precision_recall_auc, precision_score).

Run Name	Created	recall_score	rmse
powerful-gull-150	3 days ago	0.929	-
trusting-colt-811	7 days ago	1	-
judicious-elk-333	7 days ago	356	0.674
worried-wren-72	7 days ago	241	0.732
agreeable-eel-729	7 days ago	-	-
caring-croc-360	7 days ago	109	0.793

7. Compare two Runs of same model (here, Linear Regression model sample) with different alpha and l1_ratio hyperparameters

Duration:	3.1s	3.2s
▼ Parameters		
<input checked="" type="radio"/> Show diff only		
alpha		
l1_ratio		
▼ Metrics		
<input checked="" type="radio"/> Show diff only		
mae		
r2		
rmse		
➤ Tags		

Competitive Analysis

S.no	Feature	Vertex AI	MLflow
1	Python programming language support	Yes	Yes
2	R programming language support	Yes	Yes
3	Read any data source using SQL syntax	Metadata only	Uses Hive metastore
4	Ability to save, load, list, tag and version multiple models	Yes (training models and model server versions)	Yes (training models version)
5	Store metrics and details of ML model training	Yes	Yes (MLFlow Experiments)
6	Feature store - Save pre calculated tabular data to be used by other team members	Yes	No
7	Scheduling - Run notebooks, jobs or pipelines on regular intervals	Requires cloud scheduler	No
8	Integration	Natively integrated with BigQuery, Dataproc and Spark	Plugins allow to integrate with third-party storage solutions for experiments data (metrics and hyperparameters), artifacts and models
9	API Abstraction with simple one line code or User interface - Unified UI and API	Yes	Yes
10	AutoML - Try to automatically find a model that gives the best results	Yes	No

11	User Interface support for ML workflow	Supports end to end - development, serving, monitoring	Supports only for tracking and has its own registry
12	Support for open source frameworks - tensorflow, pytorch, scikit-learn	Yes	Yes
13	Use cases	<ul style="list-style-type: none"> ● Data readiness ● Feature Engineering ● Training & hyperparameter tuning ● Model serving ● Model tuning and understanding ● Model monitoring ● Model management 	<ul style="list-style-type: none"> ● Tracking ● Projects ● Models ● Registry
14	Time reduction through MLOps automation	~ 15 days of manual labor	~ 6-8 days of manual work

COST COMPARISON

Vertex AI

Training [5]

- 8 hours per model
- Per hour pricing
 - n1-standard-8, 8 CPU, 32 GB Memory = \$0.436998 / hr
 - NVIDIA_TESLA_K80 GPU = \$0.517500 / hr
 - ~\$1 per hour = ~\$8 per day per model
 - 8*365 = \$3000 per year per model

Serving [5]

Per hour pricing

- n2-standard-8 \$0.39 per hour
- 1 machine \$0.39 per hour = \$3,416 per year (if running all the time)
- Scale up to 10 machines during a hypothetical two-hour peak that occurs everyday
 - 9 machines \$0.39 per hour = \$2,562 per year (2 hrs per day)
- Total = \$6,000

Engineering / deployment

- Initial CI/CD pipeline setup cost = \$6,000

Costs for deploying 1 model to production

Model Infrastructure	Kubernetes clusters with a load balancer	\$6,000
Data Support	Independent data pipeline manager for continuous update of analytic data	\$10,000 (labor) + \$3,000
Engineering / Deployment	Continuous integration and continuous deployment (CI/CD) system to pull model from registry	\$6,000 (labor) + \$200
Total Investment / yr		\$16,000(labor) + \$9200
5 year TCO		\$62,000

Costs for deploying an additional model to production

Model Infrastructure	Cluster can be adjusted to run multiple models according to changing usage demands	+ \$2,561 / yr
Data Support	Modules can be added to pipeline management system	+ \$2,250 (labor)
Engineering / Deployment	CI/CD system can be extended to pick up additional model	+ \$750 (labor)
Total Investment / yr	\$3,000(labor) + \$2,561 / yr	
5 year TCO	+ \$15,805	

1. Cost for deploying 1 model to production = \$62,000 over first 5 years
2. Cost for deploying 2 models to production = \$77,805 over first 5 years
3. Cost for deploying 100 models to production = \$1.63 Million over first 5 years

MLFlow

Training

- Same as existing architecture

Serving

- Same as existing architecture

Engineering/deployment

- Initial CI/CD pipeline setup cost = \$18,000

Costs for deploying 1 model to production

Model Infrastructure	Kubernetes clusters with a load balancer	\$7,000
Data Support	Independent data pipeline manager for continuous update of analytic data	\$7,000 (labor) + \$3,200
Engineering / Deployment	Continuous integration and continuous deployment (CI/CD) system to pull model from registry	\$18,000 (labor) + \$300
Total Investment / yr		\$25,000(labor) + \$11,500
5 year TCO		\$81,000

Costs for deploying an additional model to production

Model Infrastructure	Cluster can be adjusted to run multiple models according to changing usage demands	+ \$3,781 / yr
Data Support	Modules can be added to pipeline management system	+ \$1,575 (labor)
Engineering / Deployment	CI/CD system can be extended to pick up additional model	+ \$2,250 (labor)
Total Investment / yr	\$3,825(labor) + \$3,781 / yr	
5 year TCO	+ \$22,730	

1. Cost for deploying 1 model to production = \$81,000 over first 5 years
2. Cost for deploying 2 models to production = \$103,730 over first 5 years
3. Cost for deploying 100 models to production = \$2.3 Million over first 5 years

RECOMMENDATION

The best alternative is **Vertex AI** because,

- Vertex AI provides a **unified platform** that can accommodate data readiness to ML model deployment. MLFlow doesn't streamline the entire ML lifecycle.
- With Vertex AI, a single Data Scientist could close the product cycle from research to production, which **saves the labor cost and time** for the entire lifecycle. MLFlow only partially streamlines the deployment process and uses the current system infrastructure for deployment. So the time taken and operational/labor cost is just a little lesser than the current system despite being open source.
- Vertex AI natively uses Google infrastructure and our system is majorly on Google Cloud. So the adoption of Vertex AI to our current system is more **compatible**.

PROCESS MODEL & PLANNING ACTIVITIES

Project Phases

This is a ‘software purchase and adoption’ project to streamline the ML model deployment. The project can be divided into 2 main phases.

1. Research
 - a. Research on ML Operations
 - b. Identify a data science specific use case
 - c. Research on the working of the current system in the context of the use case
 - d. Identify the problems with the current system
 - e. Requirements gathering
 - f. Research on 2 managed services evaluating different factors
 - g. Identify the best managed service for our project and prepare guidelines on how it can be adopted into our system.
 - h. Produce final report
2. Implementation of the adopted system (Demo)

This phase is an end-to-end Data Science project by itself. The phases included are,

1. Data understanding
2. Data preparation
3. Modeling
4. Evaluation
5. Deployment

We have 2 phases with different functionalities. So we need to identify a model that works best for both documentation heavy tasks and the implementation of a data science project by itself.

Choice of Process Model

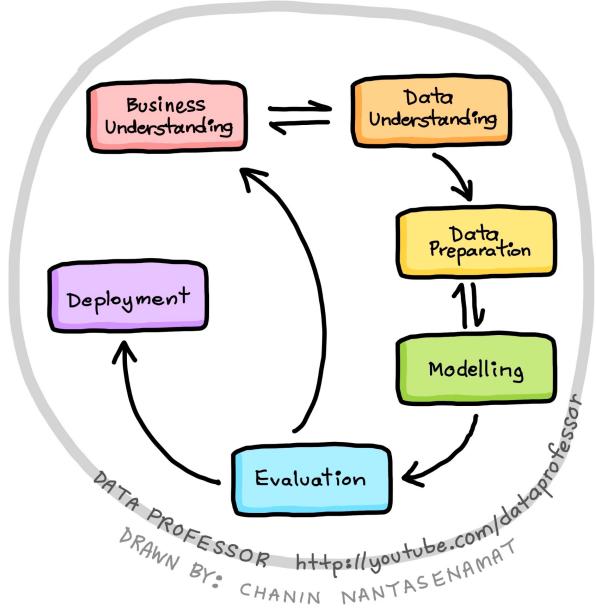


Figure 11 - CRISP-DM Model ([Image Credit](#))

CRISP-DM works best for the project since,

1. It is generalizable among Data Science projects
2. A 'mix & match' model which accommodates both horizontal and vertical slicing, i.e, waterfall and agile
3. In our project, the research phase follows a horizontal slicing which helps the team to move to the implementation phase with a clear understanding of the business needs
4. The implementation phase follows a vertical slicing which helps the team to deliver frequently and keep improving the model
5. Small team size (4 members)

PROJECT PLANNING & TRACKING

CRISP-DM Horizontal Slicing (Waterfall) [8]

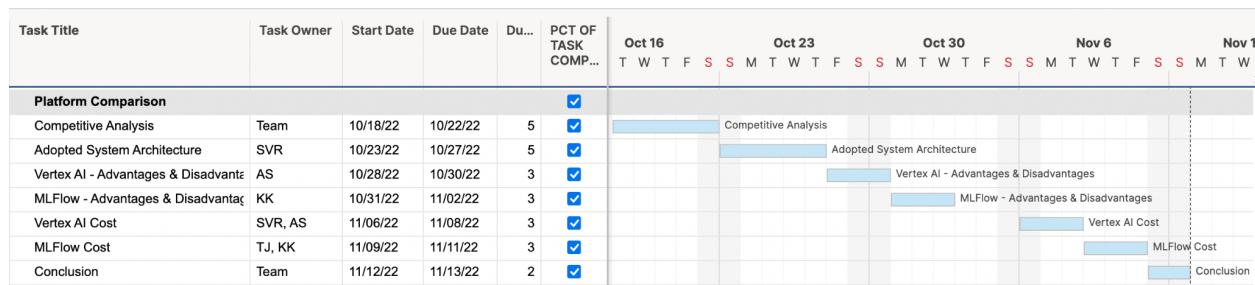
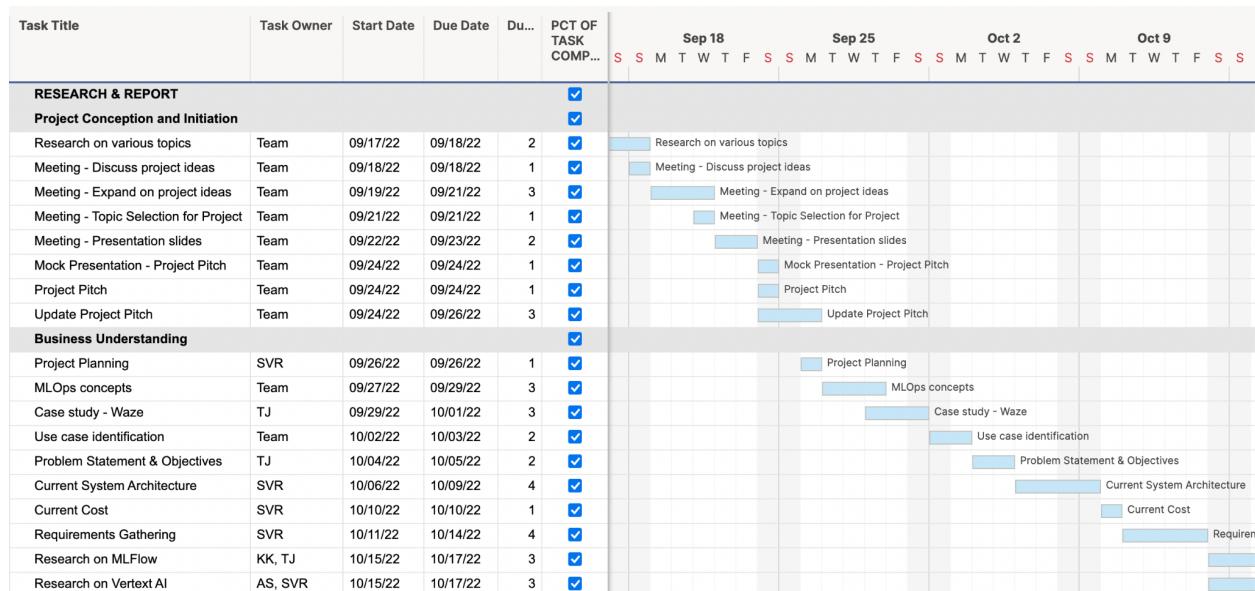
The research phase of the project includes,

1. Research on managed services and MLOps topics for the project
2. Case study of Waze
3. Research on current system architecture, cost and problems
4. Report writing

In the research phase of the project,

- Upcoming tasks are dependent on the results from the previous tasks
- Heavy emphasis on documentation was followed

Considering the above factors, a horizontal slicing(waterfall) in the CRISP-DM model works best.



Task Title	Task Owner	Start Date	Due Date	Du...	PCT OF TASK COMP...	Nov 13	Nov 20	Nov 27	Dec 4									
						M	T	W	T	F	S	S	M	T	W	T	F	S
PROJECT DOCUMENTS																		
Deliverables																		
Project Report	SVR, TJ	11/14/22	11/16/22	3	✓													
Project Summary	SVR, TJ	11/17/22	11/17/22	1	✓													
Project Presentation Slides	AS, KK	11/17/22	11/18/22	2	✓													
Team performance evaluation	SVR	11/18/22	11/18/22	1	✓													
Project Presentation																		
Mock Presentation	Team	11/18/22	11/18/22	1	✓													
Project Demo	Team	11/19/22	11/19/22	1	✓													

CRISP-DM Vertical Slicing (Agile) [8]

The implementation phase of the project includes,

1. Working demo on Vertex AI
2. Working demo on MLFlow

In this phase,

- The team had a clarity on the project requirements from the research phase.
- The team worked on the demo iteratively with heavy emphasis on construction activities.
- The demos were optimized in the future sprints
- The team had a 15 mins standup meetings daily to update the demo status
- A portion of the overall demo was expected to be implemented at the end of each week

Considering the above factors, a vertical slicing(agile) in the CRISP-DM model works best.

Task Title	Task Owner	Start Date	Due Date	Du...	PCT OF TASK COMP...	Oct 16	Oct 23	Oct 30	Nov 6									
						M	T	W	T	F	S	S	M	T	W	T	F	S
Data Understanding																		
Collect initial data	AS, KK	10/18/22	10/20/22	3	✓													
Describe data	AS, KK	10/21/22	10/21/22	1	✓													
Explore data	AS, KK	10/22/22	10/23/22	2	✓													
Verify data clarity	AS, KK	10/23/22	10/23/22	1	✓													
IMPLEMENTATION																		
Sprint 1																		
Data Preparation	AS, KK	10/24/22	10/30/22	7	✓													
Modeling	AS, KK	10/24/22	10/30/22	7	✓													
Evaluation	AS, KK	10/24/22	10/30/22	7	✓													
Deployment	AS, KK	10/24/22	10/30/22	7	✓													
Sprint 2																		
Data Preparation	AS, KK	10/31/22	11/06/22	7	✓													
Modeling	AS, KK	10/31/22	11/06/22	7	✓													
Evaluation	AS, KK	10/31/22	11/06/22	7	✓													
Deployment	AS, KK	10/31/22	11/06/22	7	✓													
Sprint 3																		
Data Preparation	AS, KK	11/07/22	11/13/22	7	✓													
Modeling	AS, KK	11/07/22	11/13/22	7	✓													
Evaluation	AS, KK	11/07/22	11/13/22	7	✓													
Deployment	AS, KK	11/07/22	11/13/22	7	✓													
Sprint 4																		
Data Preparation	AS, KK	11/14/22	11/18/22	5	✓													
Modeling	AS, KK	11/14/22	11/18/22	5	✓													
Evaluation	AS, KK	11/14/22	11/18/22	5	✓													
Deployment	AS, KK	11/14/22	11/18/22	5	✓													

REFERENCES

- [1] [Under the Hood: Real-time ETA and How Waze Knows You're on the Fastest Route](#)
- [2] [The Life Cycle of a Machine Learning Project: What Are the Stages? - neptune.ai](#)
- [3] [Estimate the Time, Cost, and Deliverables of an ML App Project](#)
- [4] [What is the Cost to Deploy and Maintain a Machine Learning Model](#)
- [5] [Vertex AI Pricing](#)
- [6] [MLFlow Documentation](#)
- [7] [A Comprehensive Comparison Between Kubeflow and MLflow](#)
- [8] [What is CRISP DM?](#)
- [9] [Vertex AI](#)