



Machine Learning Workflow

Michell S. Handaka

FOUNDER & CEO



[michellsh](#)



michell.s.handaka@glair.ai

Kevin Yauris

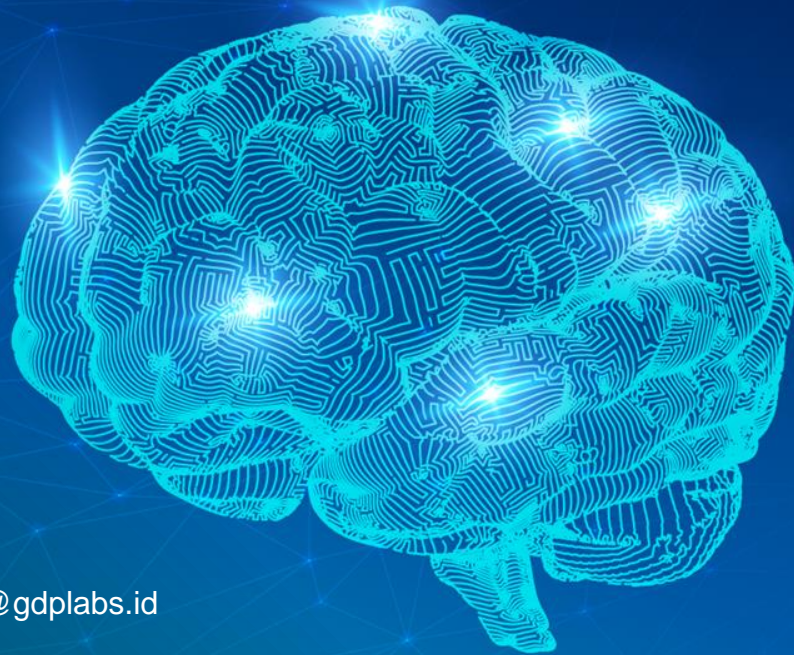
AI ENGINEER



[kevinyauris](#)



kevin.yauris@gdplabs.id



Contact Us



[glair](#)



[glair.ai](#)



hi@glair.ai



OUTLINE

- 01** Introduction to GLAIR
- 02** Launching an AI Initiative
- 03** Study Case: Propensity Modeling
- 04** Machine Learning Workflow
- 05** Q&A

Introduction to GLAIR

01

glair Products and Services

glair Consulting Service

On demand talents capable of crafting customized solutions



glair Analytics Platform

An easy & intelligent way to turn your data into insights



glair Training Center


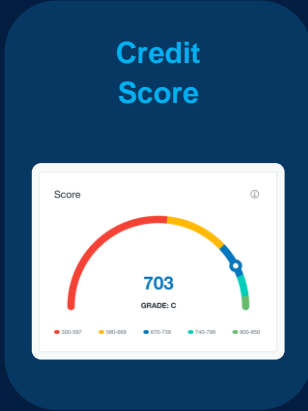
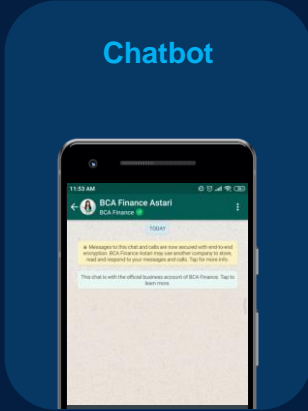
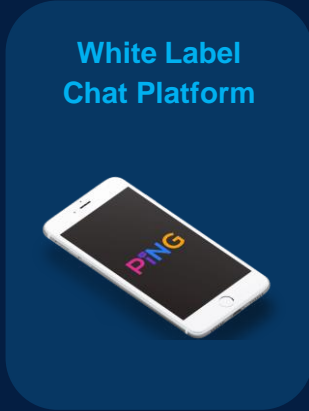
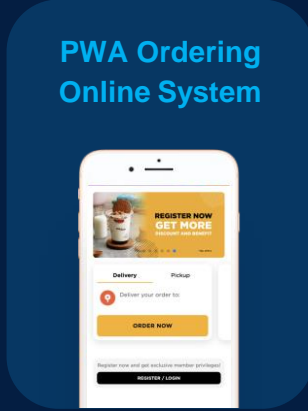
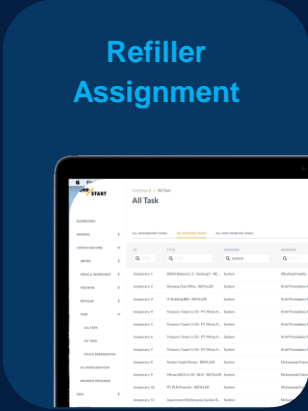
We are what we invest



The image shows a smartphone screen with a 'Fraud Detection System' interface. The top of the screen has a dark blue header with the title 'Fraud Detection System' in white. Below the header, the interface is divided into several sections. On the left is a sidebar menu with a user profile picture and name 'M. joshua' at the top. The menu items are: Dashboard, Transaction Data, Analytics, System Config, Risk Management, and User Management. The main content area on the right is titled 'Dashboard' and contains a large red box with the number '10' and the text 'High Risk Alert'. Below this is a table with columns for 'Transaction ID', 'Status', and 'Date'. The table lists several transactions, all with a status of 'High Risk' and dates ranging from 28 Aug 2024 to 02 Sep 2024.



White Label Chat Platform

A white smartphone is shown at an angle, displaying the word 'PING' in large, colorful, stylized letters on its black screen. The letters are pink, yellow, and blue. The phone is set against a dark blue background with a subtle grid pattern.[illegible]

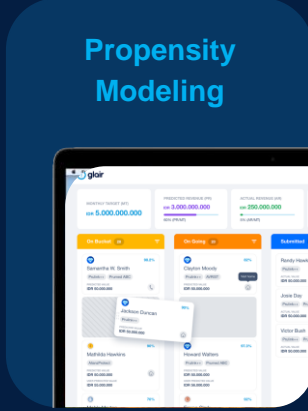
The screenshot displays the 'glor' analytics dashboard, which is used for propensity modeling. The dashboard is organized into three main sections: 'All Accounts', 'On-Site', and 'Acquisition'. Each section provides a list of accounts with their respective propensity scores and status. A modal window is open over the 'On-Site' section, showing a detailed view of a specific account's propensity score and a bar chart.

Dashboard Summary:

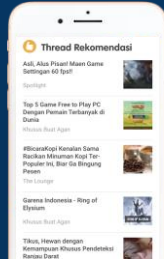
- Account Count:** 5,000,500,000
- Propensity Score Range:** 0.0000000000 to 0.0000000000
- Active Accounts:** 250,000,000

Account List (On-Site Section):

Account Name	Status	Propensity
Account A	Active	0.0000000000
Account B	Active	0.0000000000
Account C	Active	0.0000000000
Account D	Active	0.0000000000
Account E	Active	0.0000000000
Account F	Active	0.0000000000
Account G	Active	0.0000000000
Account H	Active	0.0000000000
Account I	Active	0.0000000000
Account J	Active	0.0000000000



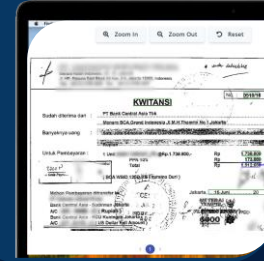
Recommendation System



OCR




Intelligent Extraction



Planogram Analytics



 **glair** Consulting Service

Sentiment Analysis & Topic Modeling



License Plate Recognition



Face Recognition



Launching an AI initiative

02

6 Preliminary Steps Before Modeling

01

Define
Objectives

02

Set
Expectations

03

Understand
the Data

04

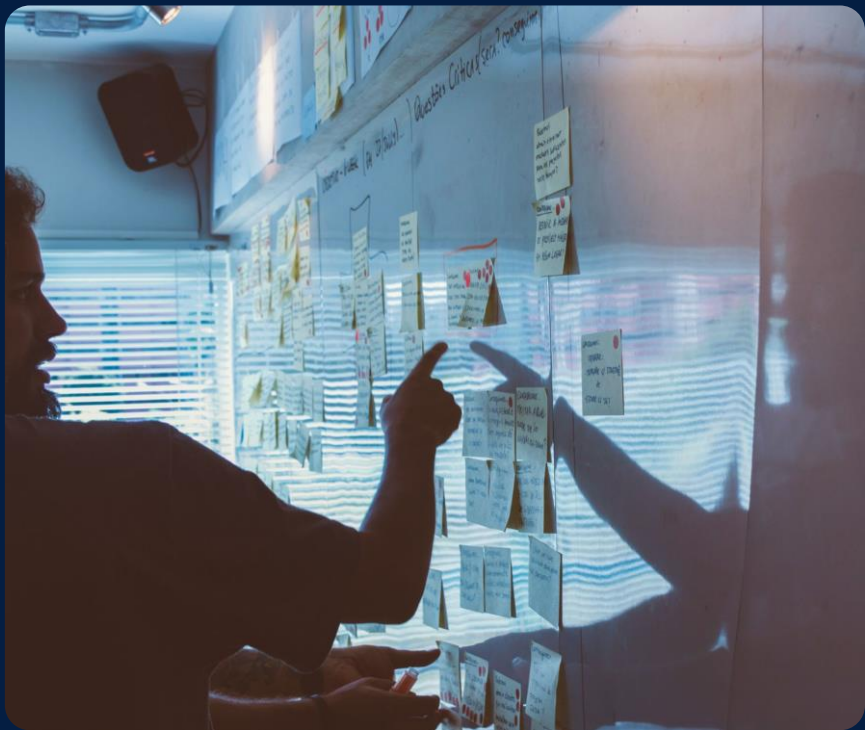
Translate the
Business Problem
Into an AI Problem

05

Determine
the Development
and Deployment
Type

06

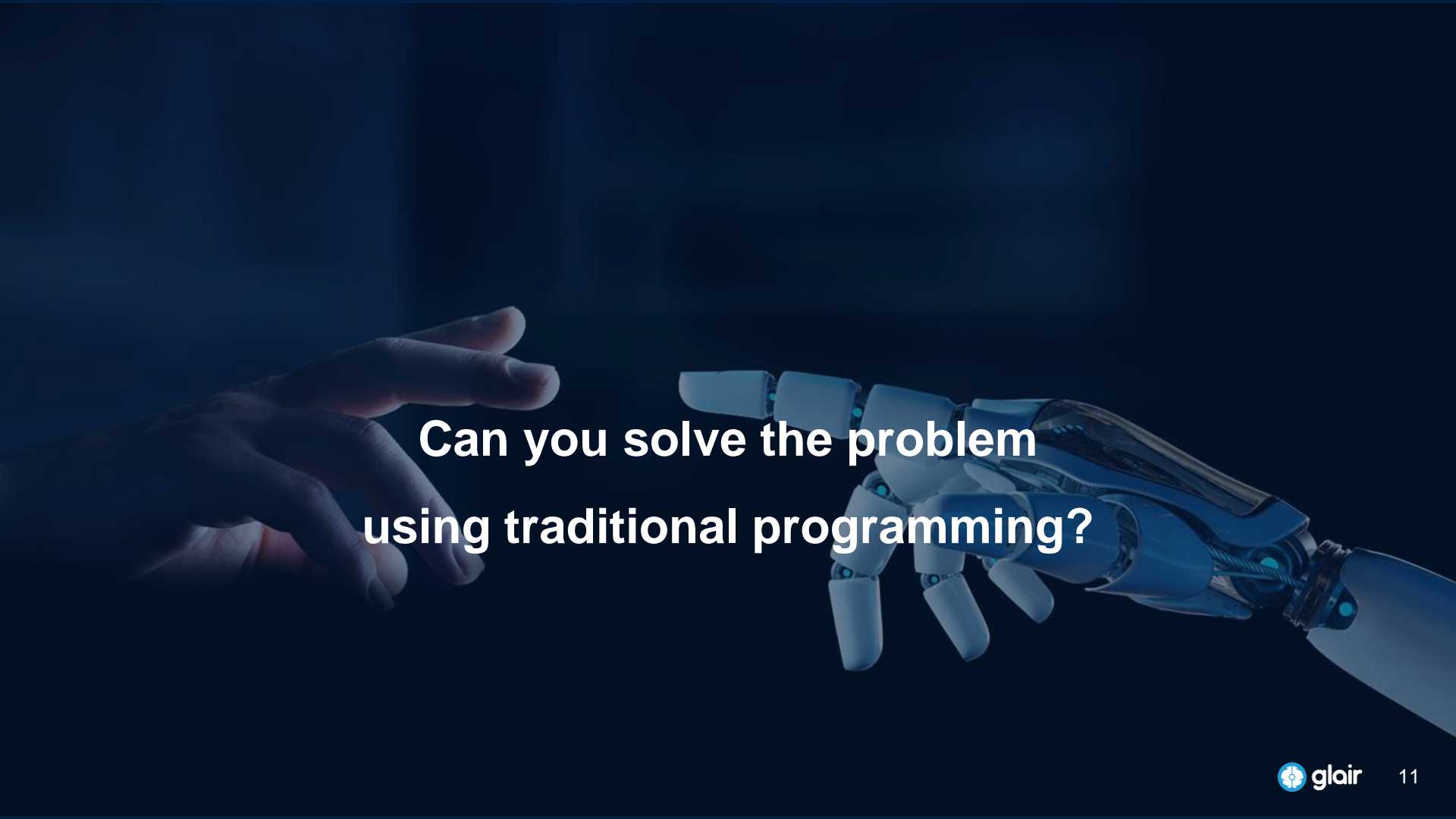
Set
Success Criteria



- Problems vs opportunities
- Increase revenue vs reduce cost
- Improvements vs innovations
- Cost and time vs differentiation

- Do you know the use cases in your industry?
- Have others approached and solved the problem using AI before?
 - What are the state of the art algorithms in the industry?
 - How about in academic setting?
 - What libraries are available?
 - What are the challenges others overcome?



A human hand on the left and a blue and white robotic hand on the right are reaching towards each other, with their fingers just inches apart. The background is a dark, textured blue.

**Can you solve the problem
using traditional programming?**



- Calculate your resources
 - Time
 - Engineers
 - Computing power

- What data do you need and what data do you have?
- Do you have enough?
 - For image classification
~1000 images for each class
 - For structured data classification
at least 100 rows, 10 rows for each class
- Where do you store the data and how do they flow?
- What is your data quality?

Data Size	Library
Big Data	 
Small Data	 

What type of data do you have?

- Structured (Tables)
- Unstructured (Image, Text, Speech)

The data type will affect the algorithm and library choices (including data processing)

Data Type	Algorithm	Library	SaaS
Structured	Logistic Regression, Decision Tree Family, GLM	Pandas, Apache Spark, Scikit-learn, XGBoost	Google AutoML Table, Microsoft Azure Machine Learning Studio
Image	CNN (Inception, ResNet, ...)	Tensorflow, Keras, OpenCV	Konvergen, Amazon Rekognition, Google Cloud Vision/AutoML, Microsoft Azure Computer Vision
Text	RNN, Transformer (BERT, GPT, ...)	Tensorflow, Keras, Spacy, Gensim	Prosa.ai, Amazon Comprehend, Google Cloud Natural Language/AutoML, Microsoft Text Analytics

Business Problem	AI Problem
Credit Scoring	Binary Classification
Propensity Model	Binary Classification
Fraud Detection	Binary Classification
Demand Prediction	Regression
Sentiment Analysis	NLP: Text Classification
Face Recognition	CompVis: Image/Face Classification

Business Problem	Business Metrics	ML Metrics
Credit Scoring	Revenue and NPL Rate	F1 Score
Propensity Model	Lift Metrics	PR-AUC
Fraud Detection	Precision and FP	PR-AUC
Demand Prediction	#Out of stock order	RMSE

- PoC vs production
- On-premise vs cloud
 - Affects library choices
(e.g. cannot use cloud SaaS on-premise!)





- Useful for measuring progress
- Go for reasonable metrics
(current state of the art)

Study Case: Propensity Modeling

for Bancassurance



03



Need to understand existing
customers segmentation



Need to match right offers
with right customers

INPUT



Customer
Demographic
Data



Product
Specification
Data



Interaction Data:
Transaction Data



Contexts Data:
Promo Events



External Data:
Mortality Table

OUTPUT



Customer
Segmentation



Lead
Generation
(Classification)



Premium
Suggestion
(Regression)

Business

Data Information

CUSTOMER & TRANSACTION DATA FOR THE LAST 18 MONTHS



Customer



Interaction



Transaction



External



Product

Demographic data of existing bancassurance customers:
Credit card, auto finance, mortgage, a “main” table that
shows customers’ attributes

SAMPLE CASE

Current month : March 2021

Training data

September 2019 - February 2021



Leads data

April 2021

For example:

If it is March 2021 and we have data for the last 18 months, then September 2019 - February 2021 will be used as our training data, and leads data will be obtained for April 2021.

Business Metrics



Lift Metric

Lift metric is a measure **comparing** the **relative performance** of a **propensity model** vs. a **random guess**

Business → Business Metrics

Lift Metric for Propensity

Sample data:

Rank	Actual	Approached?
1	1	✓
2	0	✓
3	1	✓
4	1	✓
5	0	✓
6	0	✗
7	1	✗
8	0	✗
9	0	✗
10	0	✗



Bucket size is number of customers that will be approached and determined according to the company's capabilities

The **lift metric** is defined as the ratio of model score to random guess score

1. There are 10 customers
2. The model will rank these customers based on their probability to convert (higher rank indicates higher chance to convert)
3. Say we decide to approach 5 customers (the bucket size is 5)
4. After approaching the top 5 customers, only 3 of them converted
In this case, our model's score is: $3/5 = 0.6$
5. Assuming that we approach all 10 customers (random guess), we will have 4 converting customers. So the random guess score is: $4/10 = 0.4$
6. The model's lift is then $\frac{0.6}{0.4} = 1.5$

Benefits



Automated Pipeline

to create propensity
models and make
predictions anytime



Robust

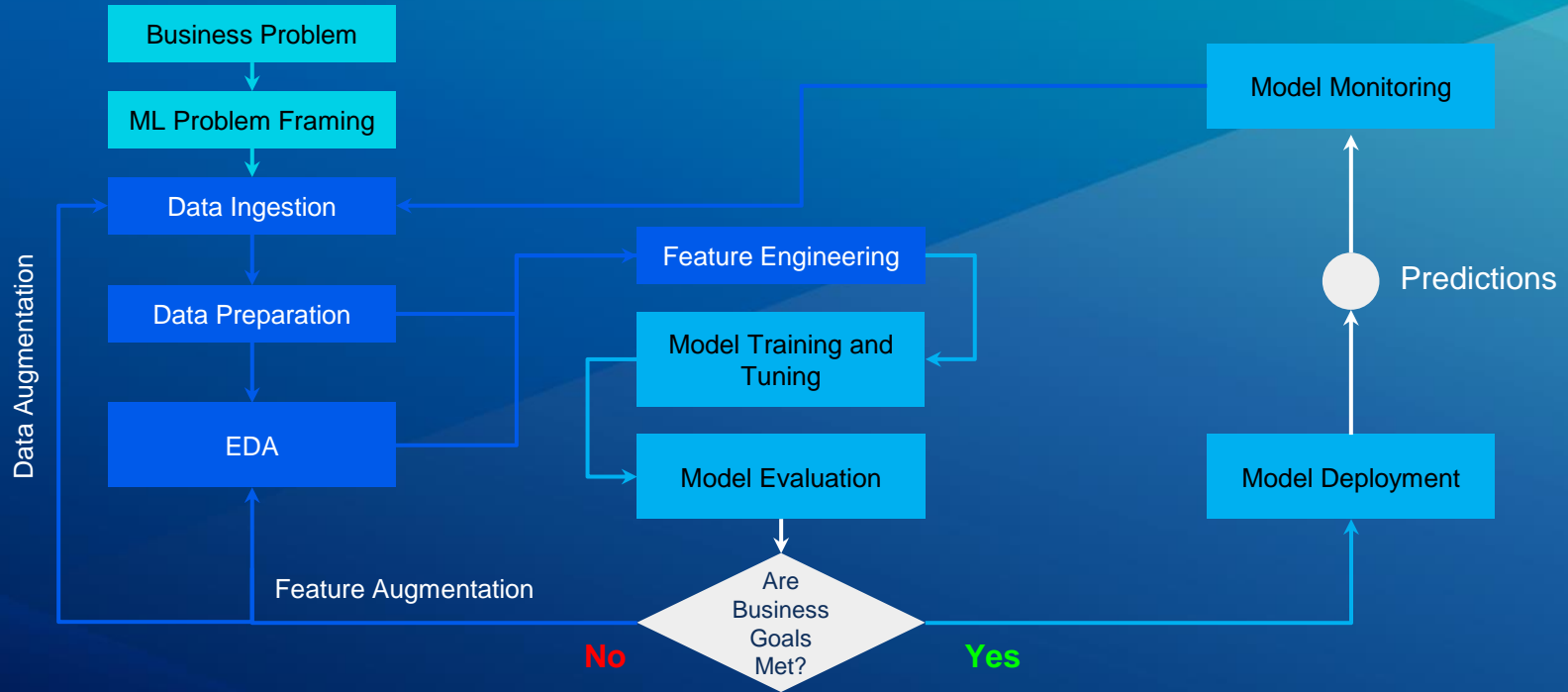
to changes in data
distribution or behavior



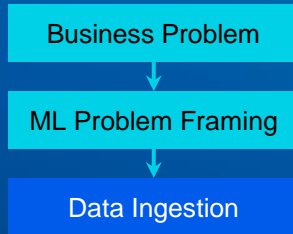
Extensible

for new data or
custom processes
to be added easily

Machine Learning Workflow



Machine Learning Workflow



- **Transportation of data from assorted sources**
to a **storage medium** where it **can be accessed, used, and analyzed**
by an organization.
- The destination is **typically a data warehouse, data mart, database,**
or a document store.
- **Sources may be almost anything.**

Batch processing

the ingestion layer **periodically collects and groups** source data and sends it to the destination system.

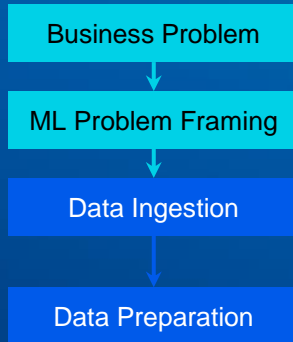
- Groups may be processed based on any **logical ordering**, the activation of **certain conditions**, or a **simple schedule**.

Real-time processing

(stream processing or streaming) involves no grouping at all.

- Data is sourced, manipulated, and loaded **as soon as it's created or recognized by the data ingestion layer**.
- This kind of ingestion is **more expensive**, since it requires systems to **constantly monitor sources and accept new information**.

Machine Learning Workflow



- Step in which the **data gets transformed, or Encoded, to bring it to such a state that now the machine can easily parse it.**

- Purpose:

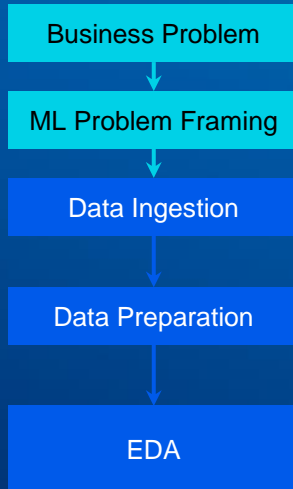
The features of the data can now be easily interpreted by the algorithm.

- Data cleansing or data cleaning:
 - **detecting and correcting** (or removing) corrupt or inaccurate records from a record set, table, or database
 - identifying **incomplete, incorrect, inaccurate or irrelevant parts** of the data and then replacing, modifying, or deleting the dirty data
- **ETL (Extract, transform, load)**

Extract, Transform, Load (ETL)

- Extract: step where **sensors wait for upstream data sources to land. Transport the data from their source locations to further transformations.**
- Transform: **apply business logic and perform actions such as filtering, grouping, and aggregation** to translate raw data into **analysis-ready** datasets.
- Load: **load the processed data and transport them to a final destination.**

Machine Learning Workflow



- Exploratory Data Analysis

Approach of analyzing data sets to summarize their main characteristics, often using **statistical graphics and other data visualization** methods.

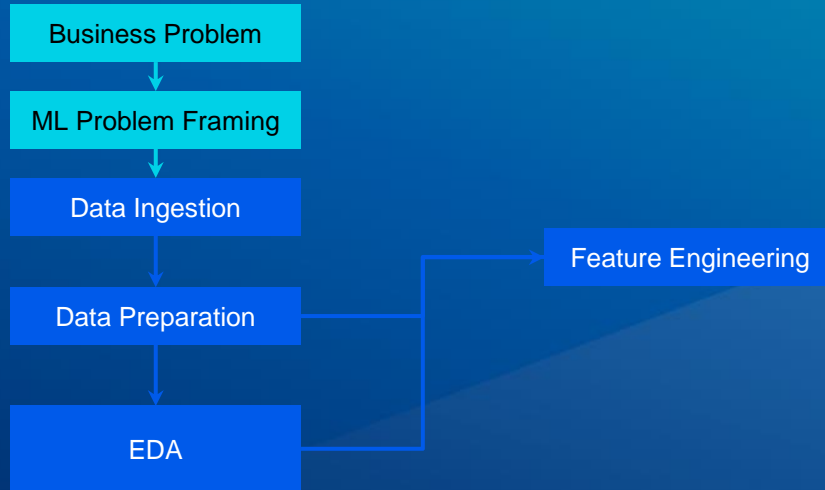
- The objectives of EDA are to:
 - **Suggest hypotheses**
 - **Assess assumptions**
 - **Support the selection of appropriate** tools and techniques
 - **Provide a basis for further data collection**

03 | Exploratory Data Analysis (EDA) Technique and Tools

- Graphical / Visualization techniques
- Dimensionality reduction
- Numerical summaries
- Statistical testing

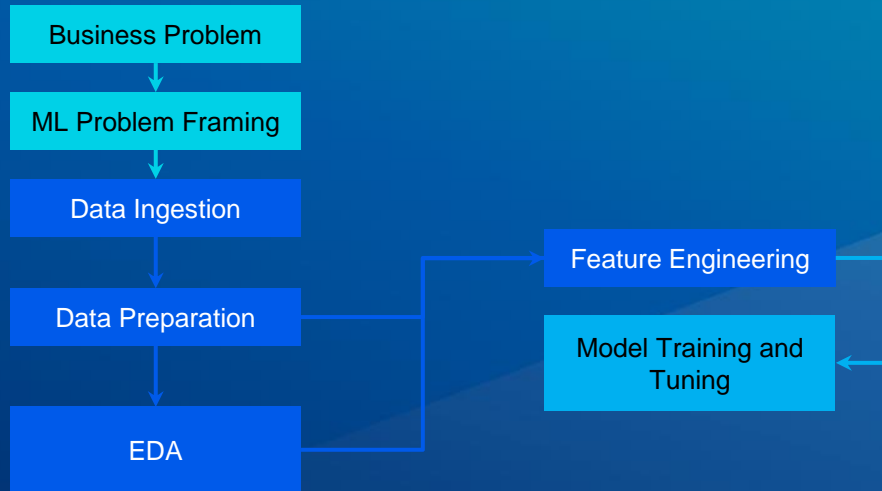


Machine Learning Workflow



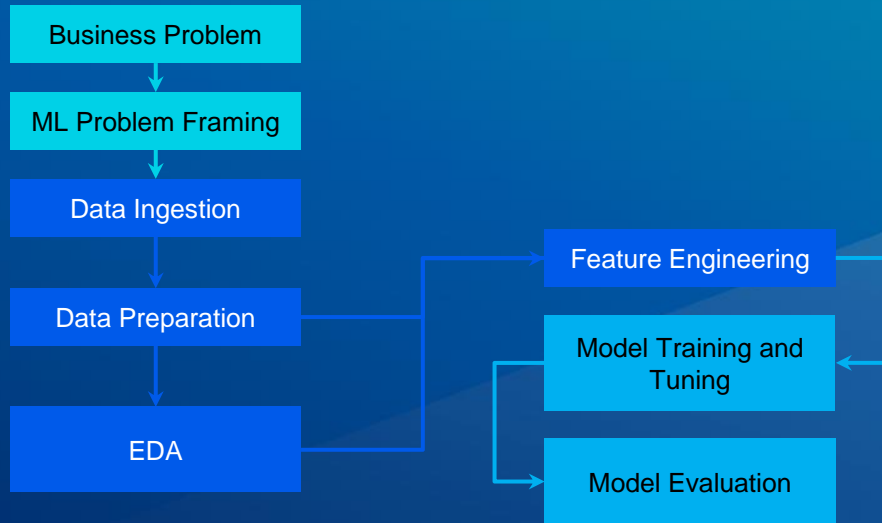
- Feature engineering is the **process of using domain knowledge to extract features** from **raw data**
- Feature engineering efforts mainly have two goals:
 - **Preparing the proper input dataset**
 - **Improving the performance** of machine learning models.
- There is a lot the of feature engineering, some of them is:
 - Imputation
 - Handling outliers
 - Binning
 - Log transformation
 - Grouping operations

Machine Learning Workflow



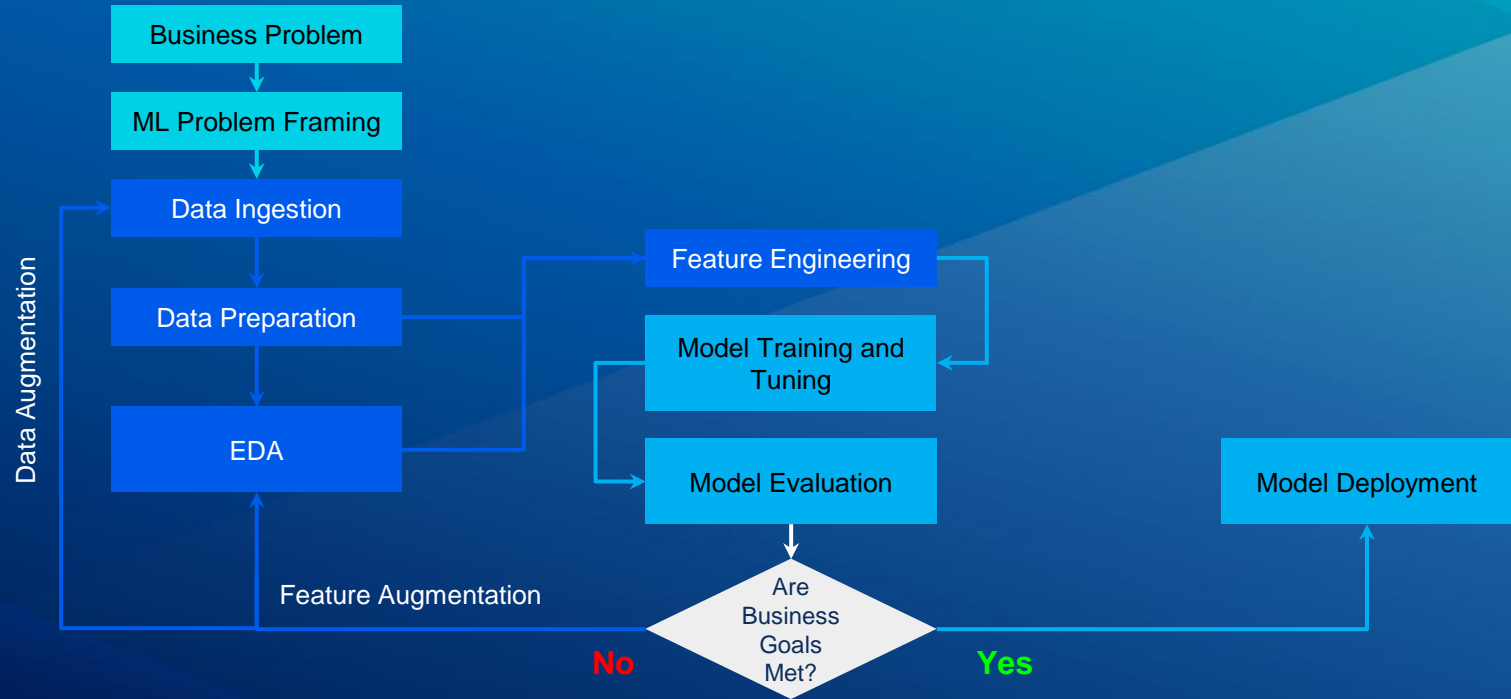
- Choose which model **to try**
- Better **try a simple model first with minimal feature engineering** and see how well it perform
- **Iteratively using a more complex model and features if needed**
- Having a **baseline model to make sure each complexity added into the model is worth to have**
- Consider **other constraints**:
 - **cost, explainability, and speed**
- **Maximize model performance by doing hyperparameter tuning**

Machine Learning Workflow



- **Evaluate** the **model** with the **test data**
- Pick a **suitable metrics** for the problem.
 - There must be a **ML and business** metrics for the problem
- If the **data is unbalanced**, **precision or accuracy is not a suitable** metrics
- If the business metrics achieved then continue with model deployment, **if not reiterate** the model creation with **data augmentation** (adding more data) or **feature augmentation** (adding other features)

Machine Learning Workflow



- **Enable the model to be used for inference**
- Considerations:
 - How to **wrap the prediction code as a production-ready service?**
 - **How to ship and load** the dumped model file?
 - Which **API / Protocol** to use?
 - **Scalability, Throughput, Latency.**
 - Deployments
 - How to **deploy new model versions?**
 - How to **rollback?**
 - Can we test it using **Canary Deployments** or **Shadow Deployments?**

ML Serving Frameworks

There are some several frameworks that provide solution



Just a REST API Wrapper

"The K8 Model Serving Projects:
KFServing and Seldon Core"



BentoML



Neptune



TensorFlow

Tensorflow Serving

cortex

Cortex



Torchserve



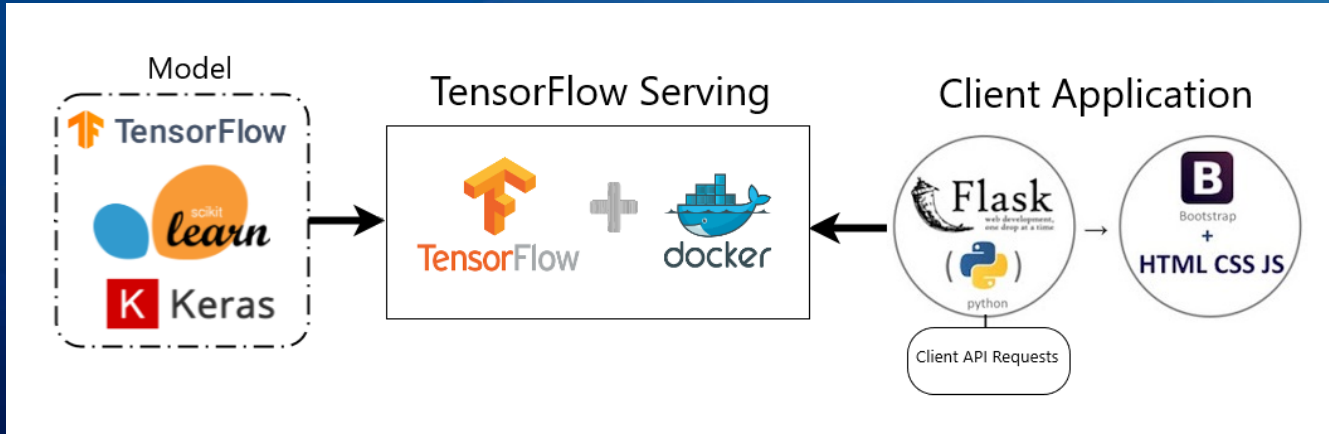
KF Serving



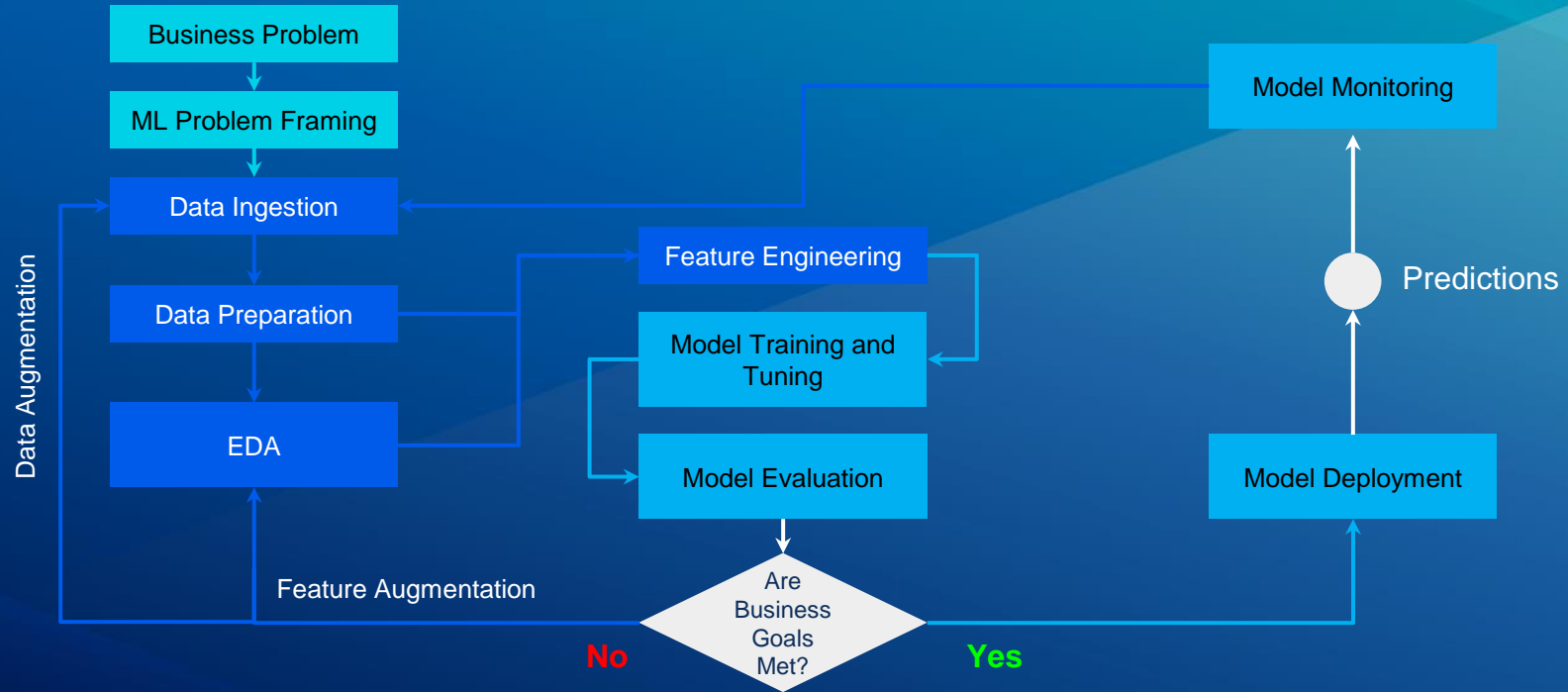
Source: Medium

TensorFlow Serving

- TensorFlow serving provides an easy integration for developers to incorporate AI in software systems and is already been used in productionizing a lot of google products.
- It can serve multiple models and multiple versions of the same model simultaneously



Machine Learning Workflow



- To make sure model still **perform well when new data coming**.
- Things to be monitored:
 - **Service Health**
 - **Data Quality & Integrity**
 - **Data & Target Drift**
 - **Bias/fairness**



Thank You! Any Questions?

Michell S. Handaka

FOUNDER & CEO



[michellsh](#)



michell.s.handaka@glair.ai

Kevin Yauris

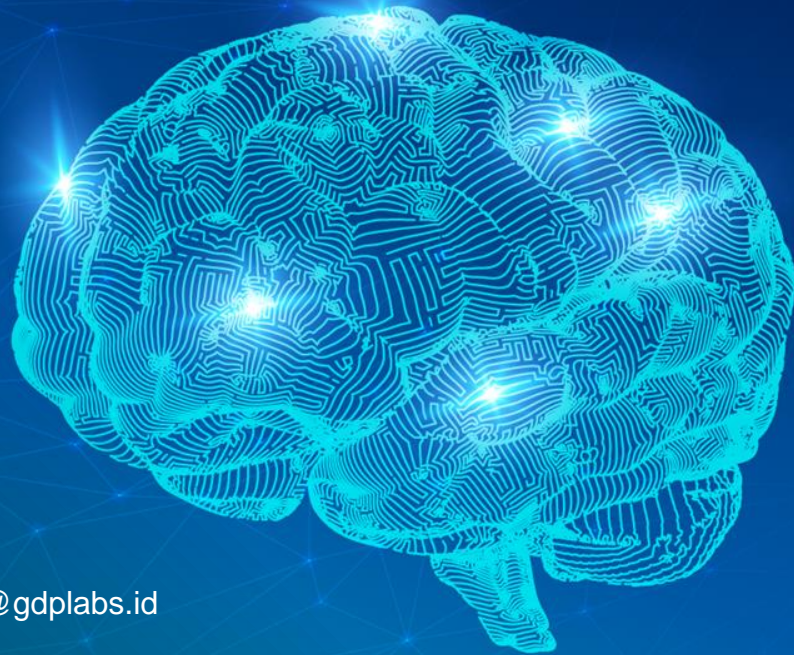
AI ENGINEER



[kevinyauris](#)



kevin.yauris@gdplabs.id



Contact Us



[glair](#)



[glair.ai](#)



hi@glair.ai