

# TensorFlow in the Real World: Deploying Models

**Gusti Triandi Winata**

Community Researcher of RL & LLM  
Cohere for AI

## About Me



## Gusti Triandi Winata

---

### Latest Work Experiences:

- Researcher, Cohere for AI 2024 - present
- MLE Consultant, Freeport Indonesia 2023 - 2024
- Mid. MLE, eFishery 2022 - 2023
- OSS Fellowship with Adobe, Major League Hacking 2021

### Education:

- Bandung Institute of Technology Graduated at 2021  
*Bachelor of Electrical and Computer Engineering*  
*Was a Bangkit Graduate of the first batch (2020)!*

# Ground Rules

Observe the following rules to ensure a supportive, inclusive, and engaging classes



Give full attention  
in class



Mute your microphone  
when you're not talking



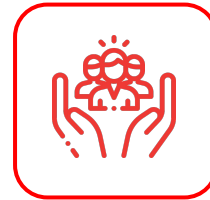
Keep your  
camera on



Turn on the CC Feature  
on Meet

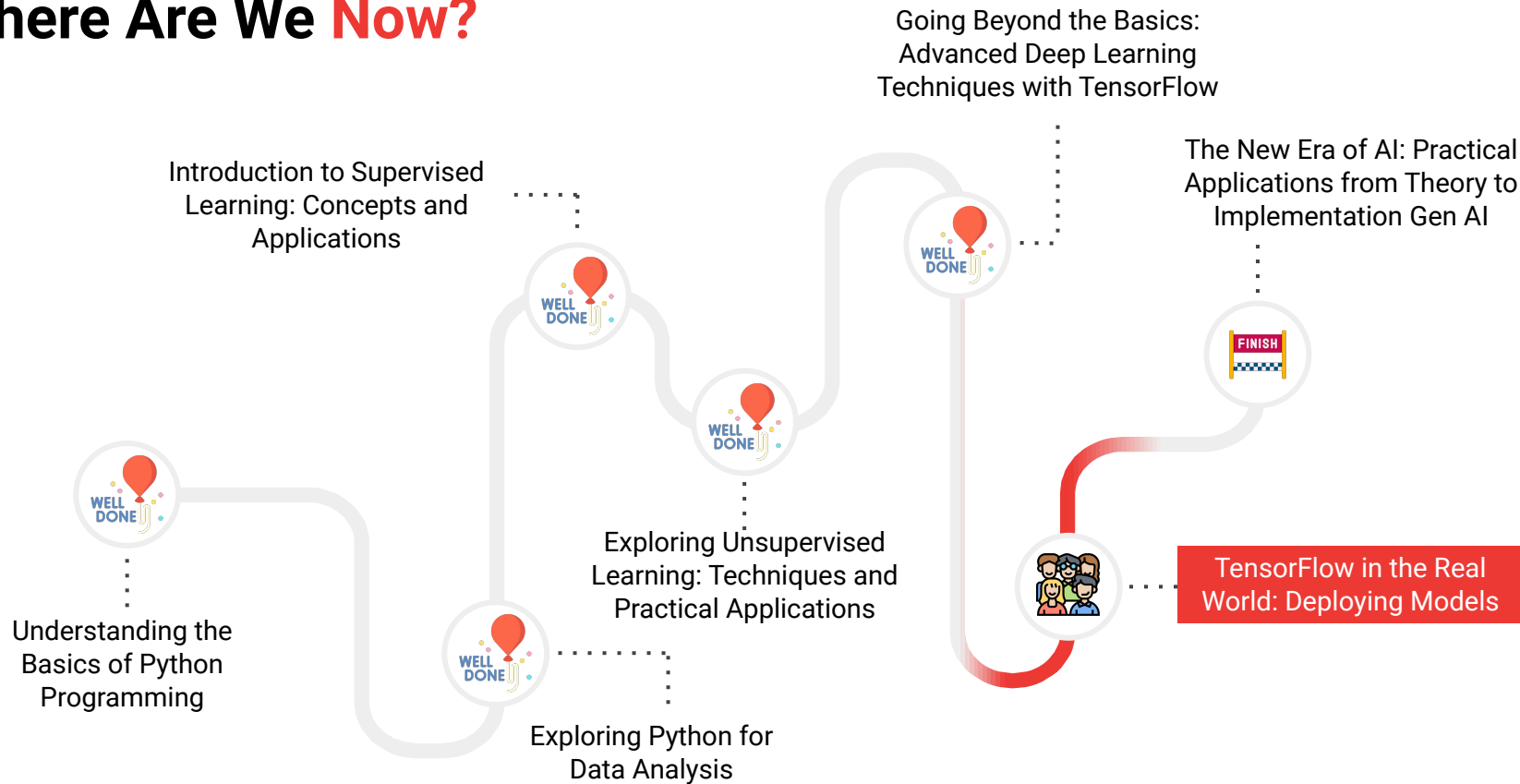


Use raise hand or chat  
to ask questions

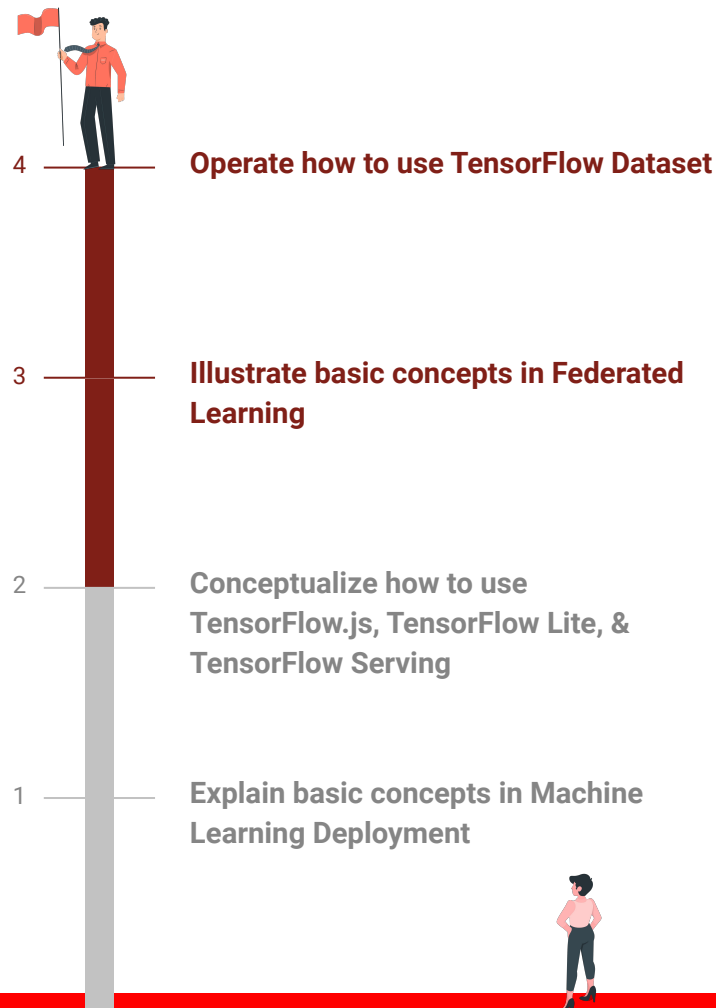


Make this room a safe place  
to learn and share

# Where Are We **Now?**



# Learning Objectives



# Today's Agenda

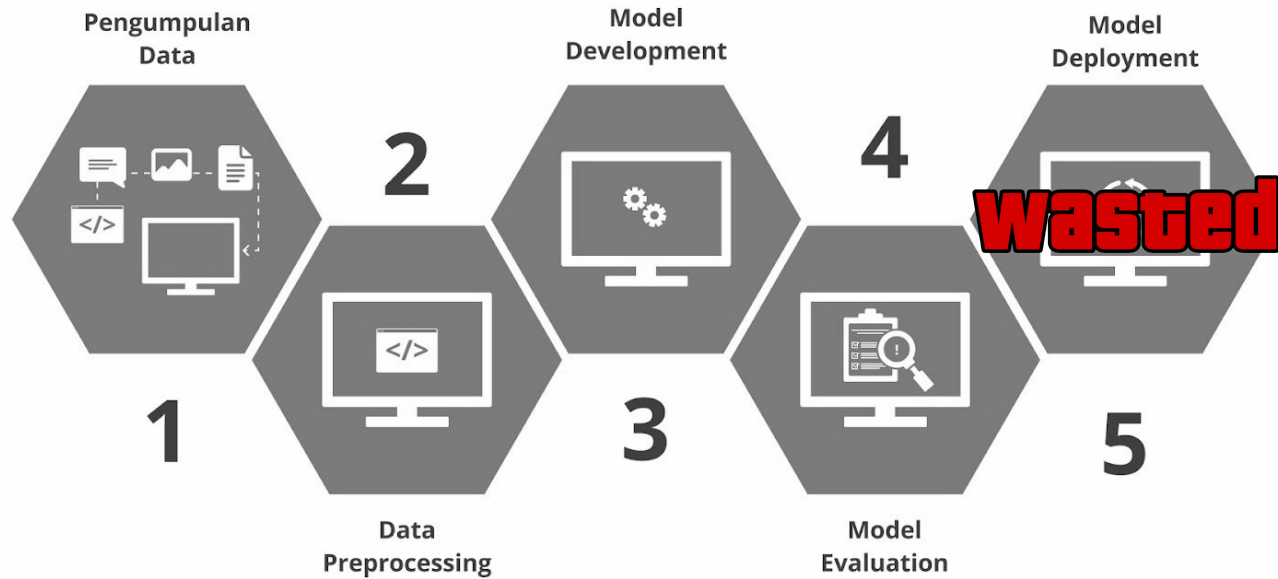
Embarking on the Journey of Machine Learning Model Deployment

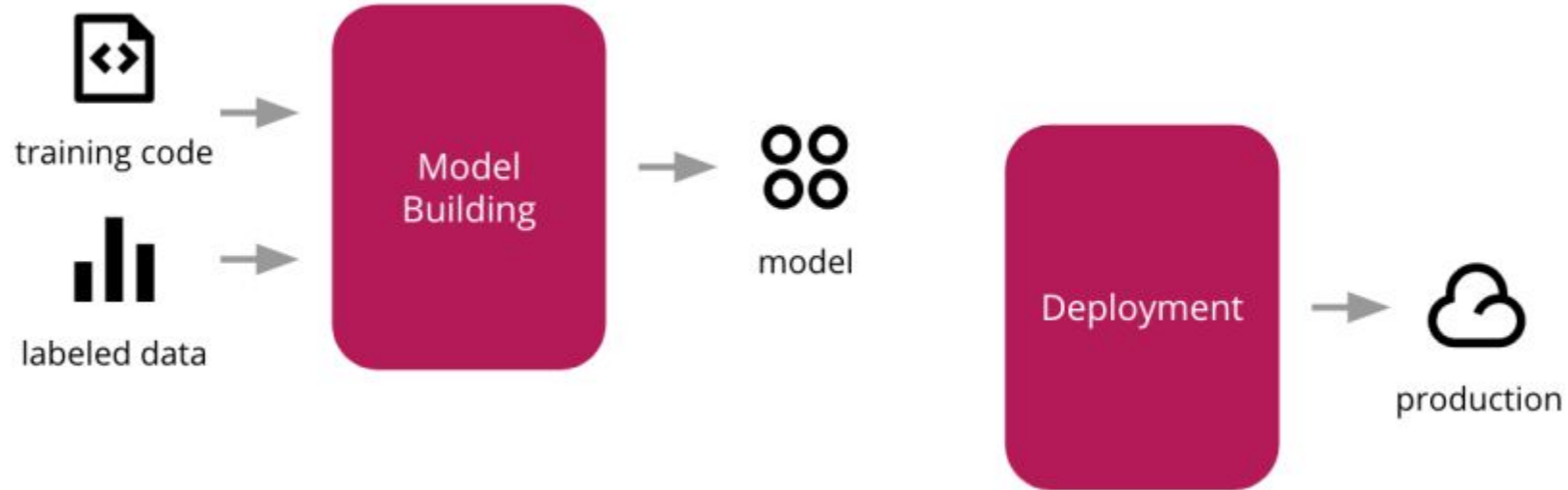
Utilizing TensorFlow Tools for Deploying Models in Real-World Applications

Highlighting Basic Principles of Federated Learning

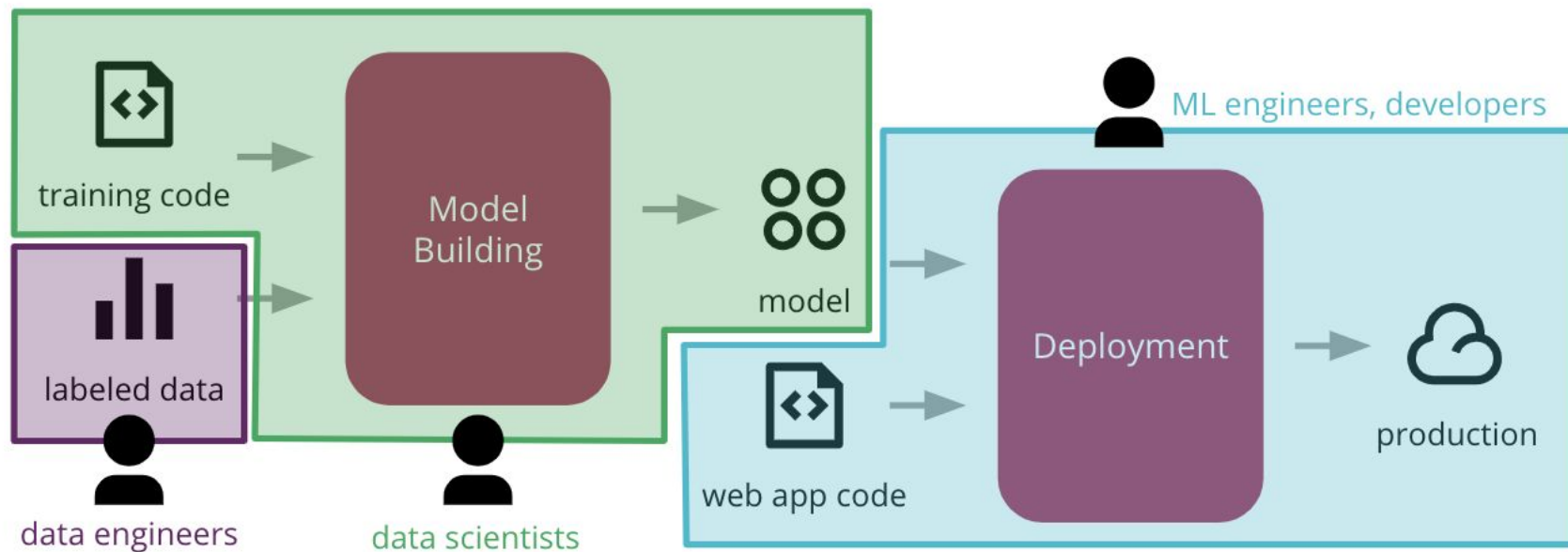
Introduction to Data Handling with TensorFlow Datasets









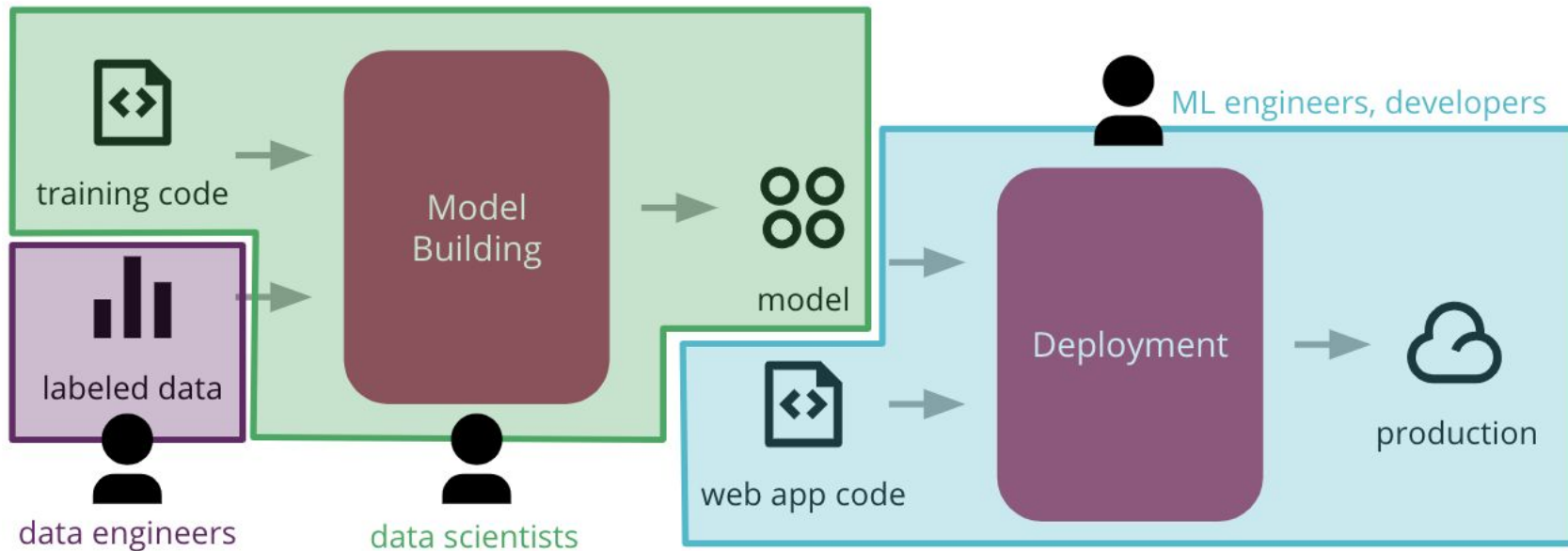


**Bangkit ML Students**



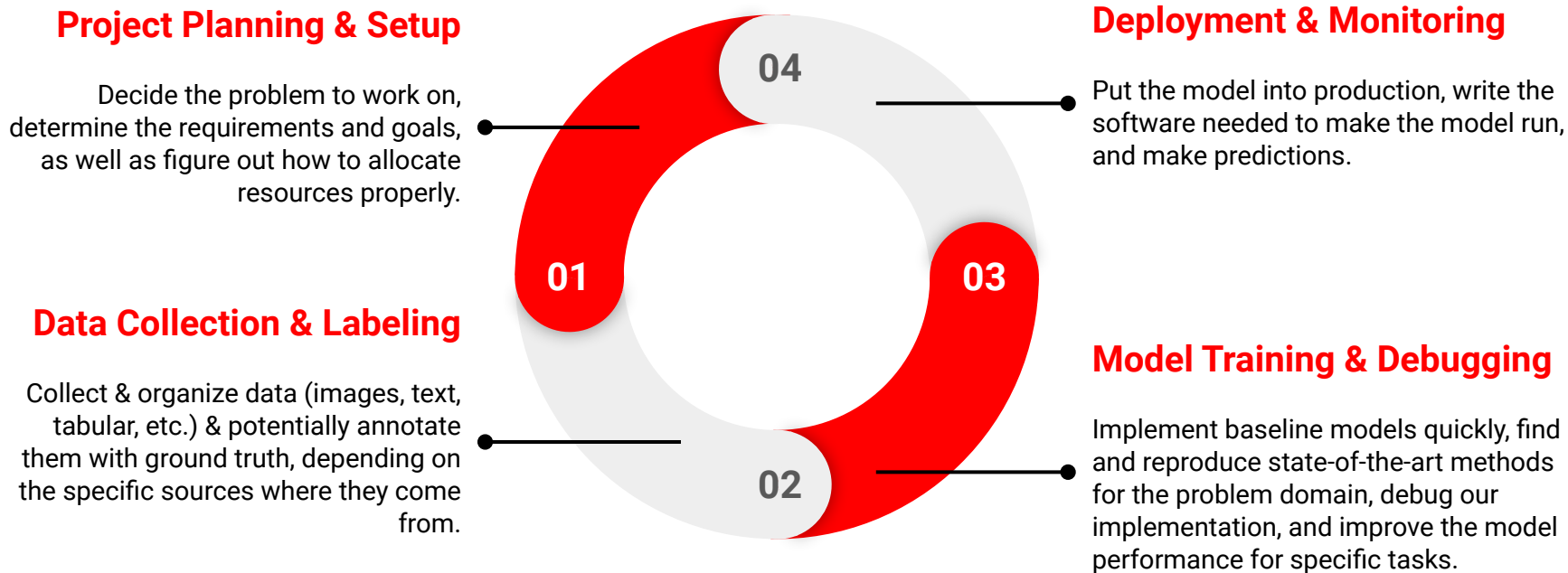


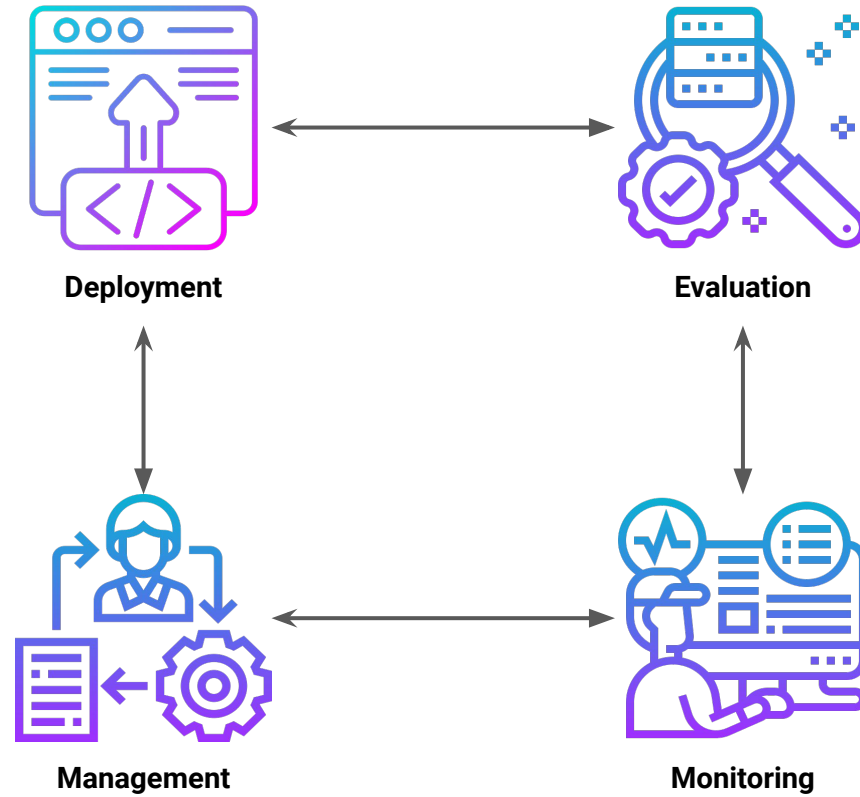
## Bangkit ML Students



# Introduction to Machine Learning Deployment

# Lifecycle of an ML Project





# The **Challenges** of Model Deployment

- ML models are sensitive to incoming data's **semantics, quantity, and completeness**.
- The performance of ML models in production **degrades over time** due to changes in actual data.
- ML models only work with data from a **specific demographic**.



# Model Deployment Options



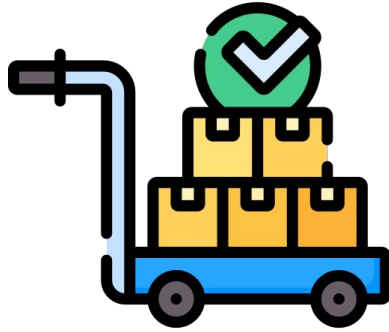
**Centralize** model in  
server



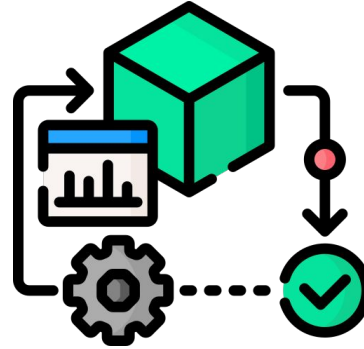
**Distribute** model on  
user device



# Model Deployment Options



Batch Prediction

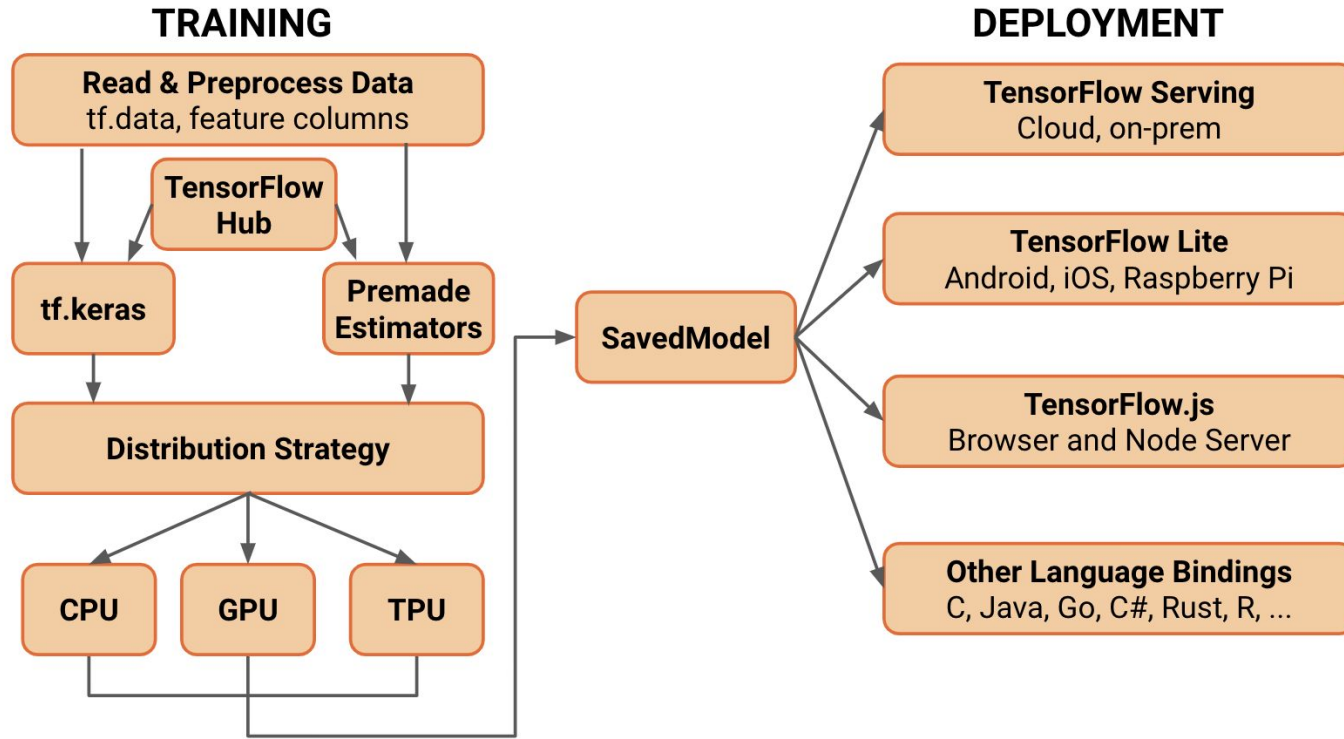


Real Time Prediction

How does the **user use** the model?



# Model Deployment Options



# Model Deployment: TensorFlow.js

# What is TensorFlow.js?

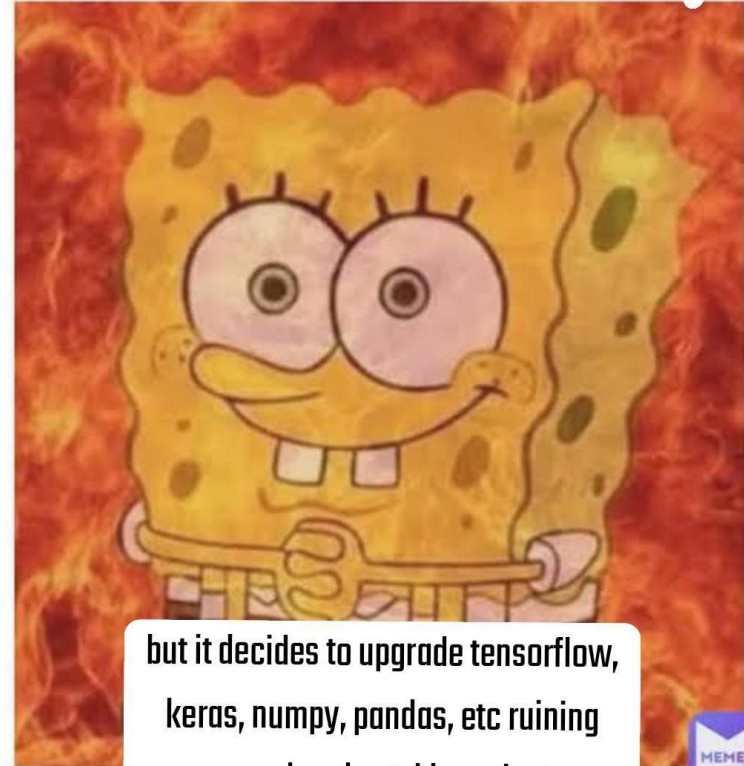
An **open-source JavaScript library** for training and deploying machine learning models in the **client's browser** or **Node.js server**.

- Run existing models
- Retrain existing models
- Develop ML with JavaScript

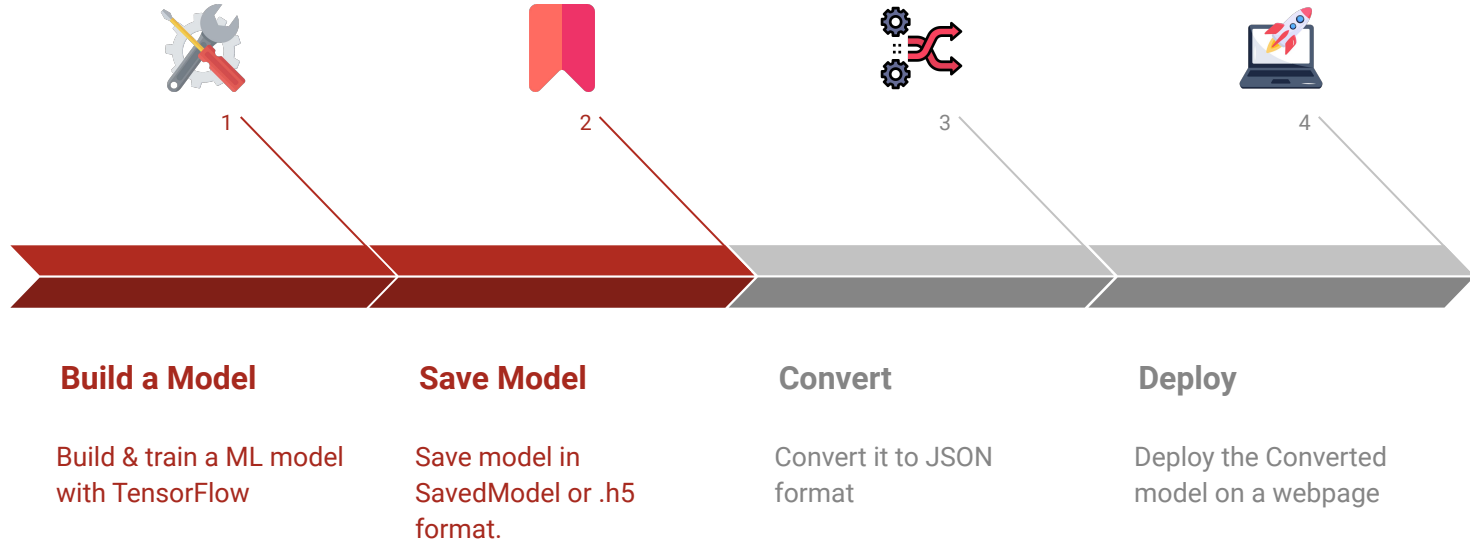


when you try to install  
a new library in python

## TensorFlow.js Disadvantages



# General Steps in TensorFlow.js



# [Hands-on] Bridging Machine Learning with Web Development



# Model Deployment: TensorFlow Lite

# What is TensorFlow Lite?

An open-source deep learning framework to **run TensorFlow models on-device**

- Optimized for on-device machine learning
- Multiple platform support
- Diverse language support
- Hardware acceleration & model optimization





# The **Challenges** of on-Device Models

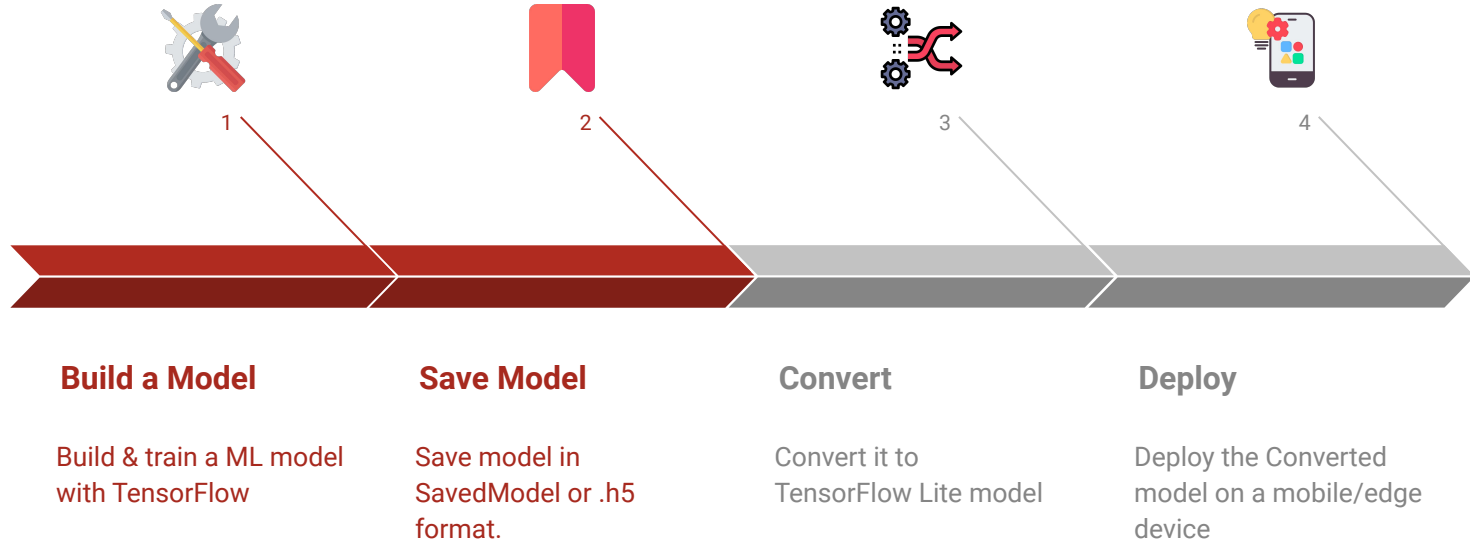
Running on-device models means we need to handle **diverse devices** & we do **not have access** to actual data.

Need to **balance** between:

- Power efficiency
- Inference latency
- Model accuracy & complexity



# General Step in TensorFlow Lite



# **[Hands-on] Streamlining TensorFlow Lite Deployment for Mobile Applications**

# **Drive Video Demonstration**



# Break Time



# Model Deployment: TensorFlow Serving

## What is TensorFlow Serving?

TensorFlow Serving is a **flexible & high-performance** serving system for machine learning models, designed for **production environments**



TensorFlow Extended

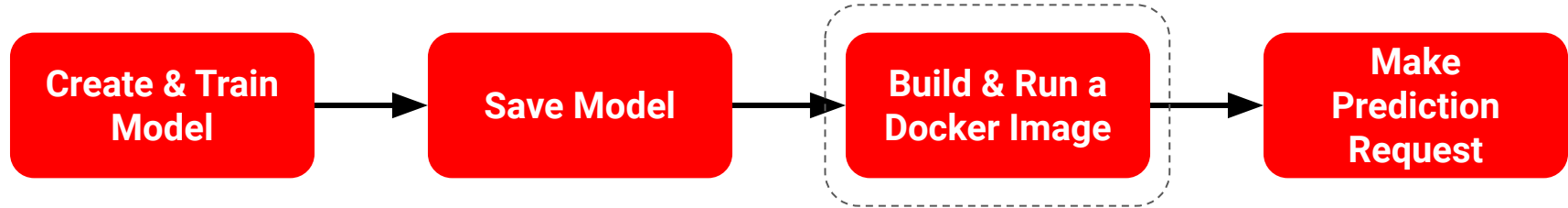
# TensorFlow Serving

TensorFlow Serving allows us to have a **centralized model**

- Easy to manage model version
- Easy to manage hardware resources based on demand
- Can have multiple serving processes



# Deploy Model with **TF Serving Workflow**



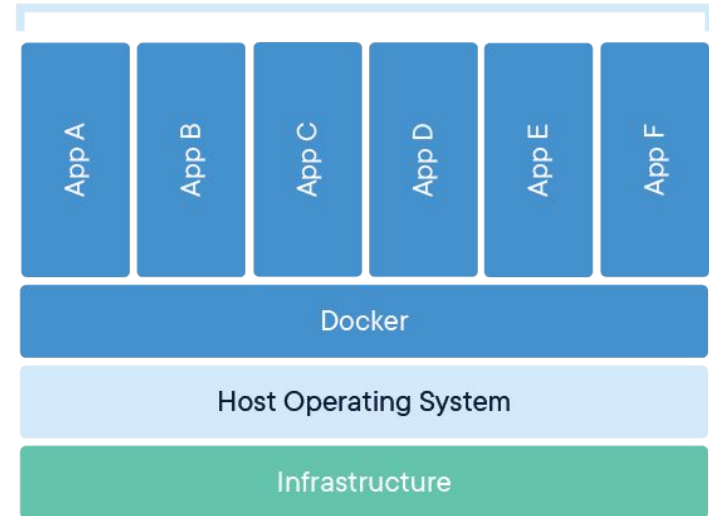
# What is **Container?** Docker?

Container is a standard unit of software that **packages up code** and **all its dependencies** so the application runs portably.

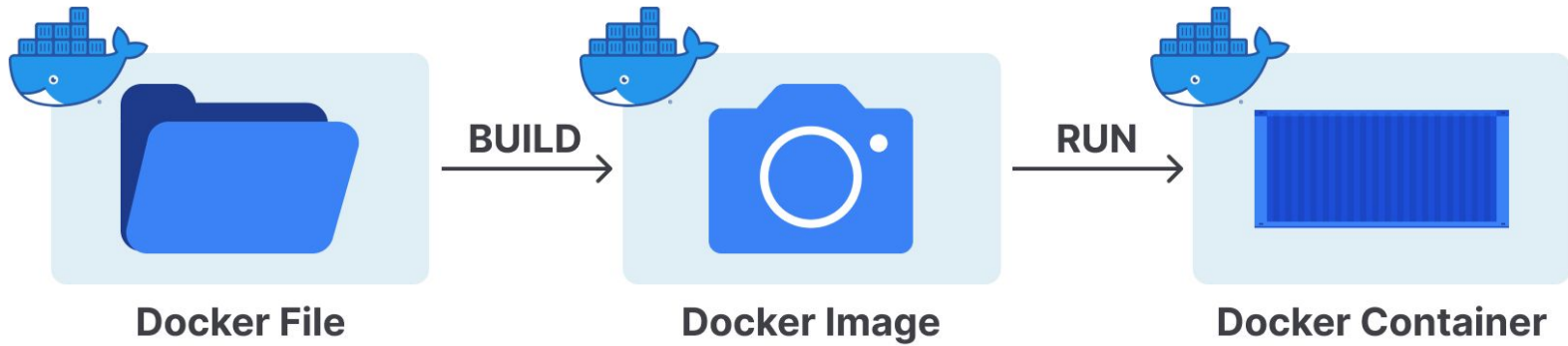
- Portable
- Lightweight
- Isolation



Containerized Applications **docker**



# How to Create **Container**?



# **[Hands-on]** **TensorFlow Serving** **Deployment in** **Cloud-native** **Environments**

# Federated Learning



# What is **Federated Learning**?

Federated learning allows each client **independently train** its own model using its own data right on the device.

Lower latency

Less power consumption

Ensuring privacy



# Deployment Option **Summary**

	TF JS	TF Lite	TF Serving	Federated Learning
<b>Model runtime</b>	Node.JS server / client's browser	On-device	Server	On-device
<b>Computing power</b>	Depends on usage	Low	High	Low
<b>Latency</b>	Depends on the model complexity	Low	Depends on the infrastructure and model complexity	Low
<b>Model complexity</b>	Depends on usage	Lighter	Heavier	Lighter
<b>Need server connection</b>	Anytime	One time only / once in a while update	Anytime	One time only / once in a while update
<b>Privacy</b>	Depends on model runtime	No need to send data to the server	Need to send data to the server	Ensuring user privacy

# Introduction to TensorFlow Datasets

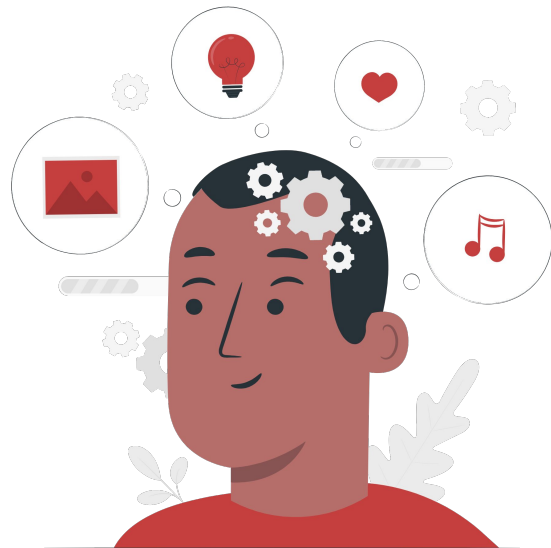
# What is TensorFlow Datasets?

TensorFlow Datasets is a **collection of datasets ready to use**, with TensorFlow or other Python ML frameworks.

**Simplicity & Performance**

**Determinism/Reproducibility**

**Customizability**



# **[Hands-on] Integrating TensorFlow Datasets into Your Machine Learning Workflow**

# Quiz

# Discussions

# Conclusion



You have learned **all the options to deploy a machine learning model** using TensorFlow.

In the end, you also learn about **TFDS** and how to use it to build an **efficient data pipeline.**



**Thank You**