# ILT 5 - TensorFlow Data and Deployment

bangkit

# Ground **Rules**

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



**Give full attention
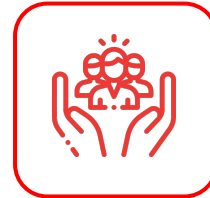in class**

**Mute your microphone
when you're not talking**

**Keep your
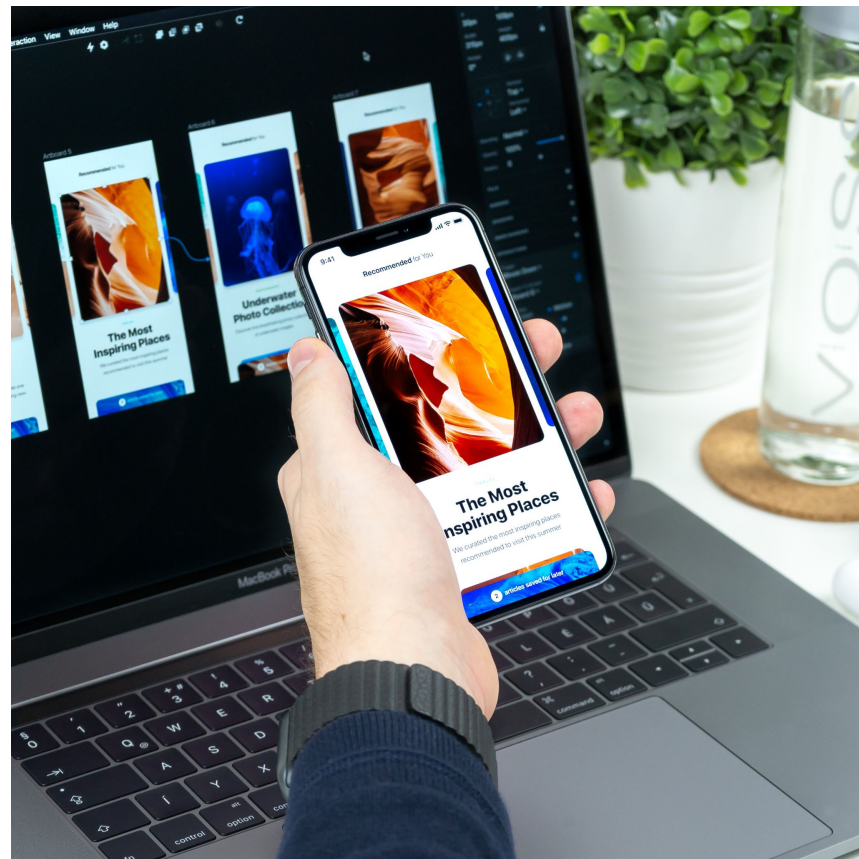camera on**

**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**

bangk!t

# Outline **Session**

- Introduction to Machine Learning **Deployment**

- Data **Pipelines**

- **Federated Learning**



**bangkit**

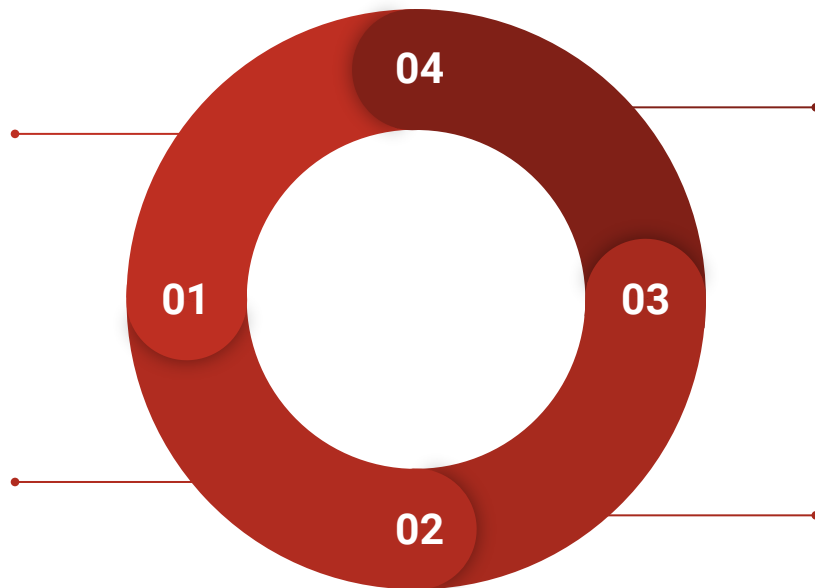# Introduction to
# ML Deployment

bangk!t

# Lifecycle of a ML Project



## Project Planning & Setup

At this phase, we want to decide the problem to work on, determine the requirements and goals, as well as figure out how to allocate resources properly

## Data Collection & Labeling

At this phase, we want to collect & organize data (images, text, tabular, etc.) & potentially annotate them with ground truth, depending on the specific sources where they come from

## Deployment & Monitoring

At this phase, we put the model into production, write the software needed to make the model run, and make predictions. We need also to monitor and maintain the system. If the data changes, we need to update the model

## Model Training & Debugging

At this phase, we want to implement baseline models quickly, find and reproduce state-of-the-art methods for the problem domain, debug our implementation, and improve the model performance for specific tasks

**bangkit**

# The **Challenges** of Model Deployment

- ML models are sensitive to the semantics, quantity, and completeness of incoming data
- The performance of ML models in production degrades over time due to changes in real data
- ML models only work with data from specific demographic



**bangkit**

# Model Deployment Options



**Centralize model in server**

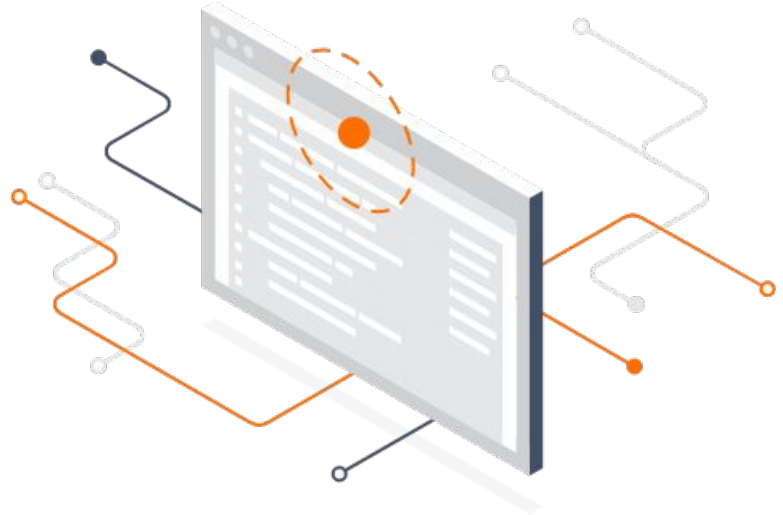**Distribute model on user device**

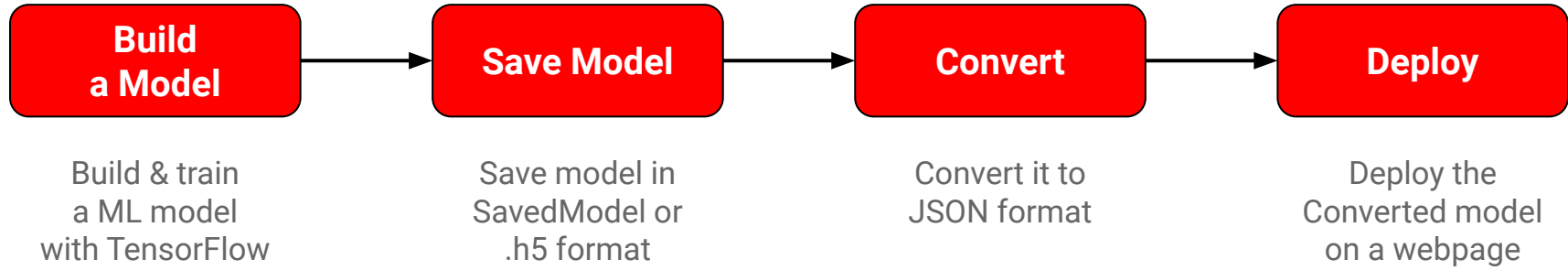bangkit

# ML Deployment
# TensorFlow.js

bangkit

# 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



bangkit

# General Steps in **TensorFlow.js**

| Build a Model | Save Model | Convert | Deploy |
|:---:|:---:|:---:|:---:|

Build a Model → Save Model → Convert → Deploy

**Build a Model**

Build & train a ML model with TensorFlow

**Save Model**

Save model in SavedModel or .h5 format

**Convert**

Convert it to JSON format

**Deploy**

Deploy the Converted model on a webpage

bangk!t

# Convert Models into JSON Format



```
model.save(saved_model_path)
```

```
!tensorflowjs_converter
--input_format=keras
saved_model_path tfjs_model_path
```

bangk!t

# What **Can** TensorFlow.js Do?

* Teachable Machine Project

# ML Deployment
## TensorFlow Lite

bangk!t

# 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
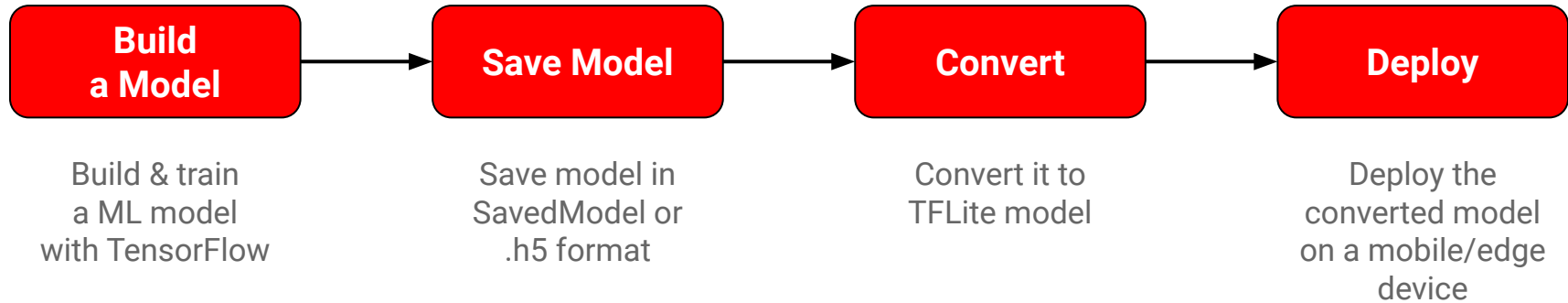


bangkit

# 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

bangk!t

# General Step in **TensorFlow Lite**

| **Build a Model** | → | **Save Model** | → | **Convert** | → | **Deploy** |
|---|---|---|---|---|---|---|

Build & train a ML model with TensorFlow

Save model in SavedModel or .h5 format

Convert it to TFLite model

Deploy the converted model on a mobile/edge device

bangk!t

# Convert Models into TF Lite Model



| TensorFlow | Saved Model | TF Lite Converter | TF Lite Model |

```python
tf.saved_model.save(model,
saved_model_path)
```

```python
converter =
tf.lite.TFLiteConverter.from_saved_mod
el(export_dir)
tflite_model = converter.convert()
tflite_model_file =
pathlib.Path('model.tflite')
tflite_model_file.write_bytes(tflite_m
odel)
```

bangk!t

# What Can TensorFlow Lite Do?

# ML Deployment
## TensorFlow Serving

bangk!t

# What is **TensorFlow Serving**?

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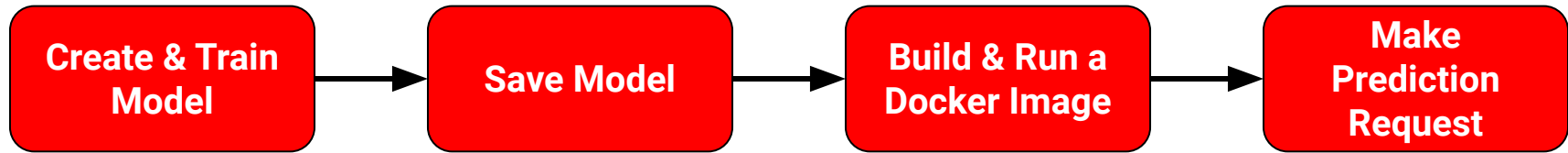


bangkit

# TensorFlow Serving

TensorFlow Serving allows us to have a centralized model

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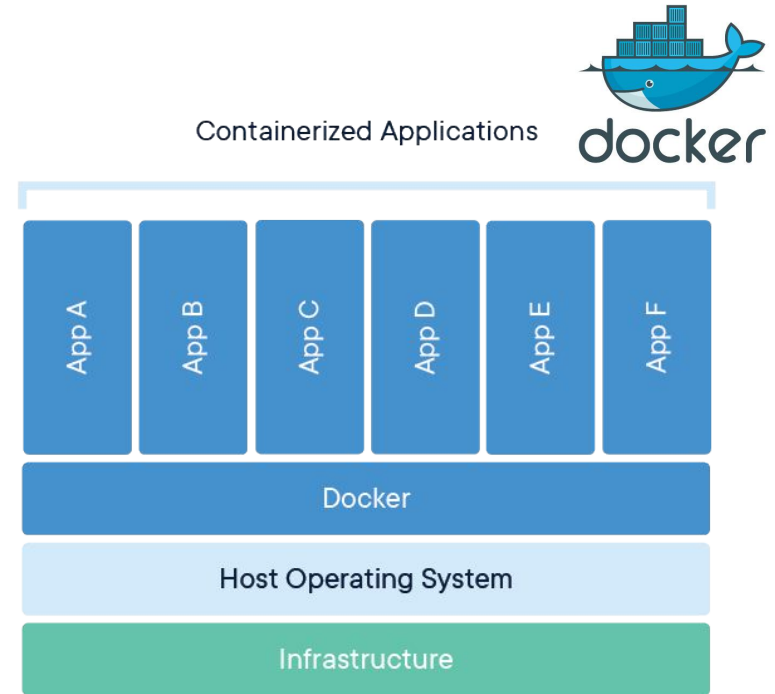- Can have multiple serving processes

bangk!t

# Deploy Model with TF Serving

| Create & Train Model | → | Save Model | → | Build & Run a Docker Image | → | Make Prediction Request |
|---|---|---|---|---|---|---|

**bangk!t**

# Whats 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

| App A | App B | App C | App D | App E | App F |

Docker

Host Operating System

Infrastructure

[Install Docker Desktop on Windows](#)

bangk!t

# How to Build & Run a Docker Image?

### Dockerfile

```
FROM tensorflow/serving:latest
```
→ Pull TF Serving image

```
COPY . /models
```
→ Copy the current directory

```
ENV MODEL_NAME=fashion-mnist
```
→ Define the MODEL_NAME

### Run the Docker image

```
docker build -t
fashion-mnist-tf-serving .

docker run -p 8080:8501
fashion-mnist-tf-serving
```

bangk!t

# How to Make Prediction Request?

This example demonstrates how to make prediction requests.

```python
import json
import requests
json_data = json.dumps({"instances": image.tolist()})


endpoint = "http://localhost:8080/v1/models/fashion-mnist:predict"


response = requests.post(endpoint, data=json_data)
prediction = tf.argmax(response.json()["predictions"][0]).numpy()
print(prediction)
```

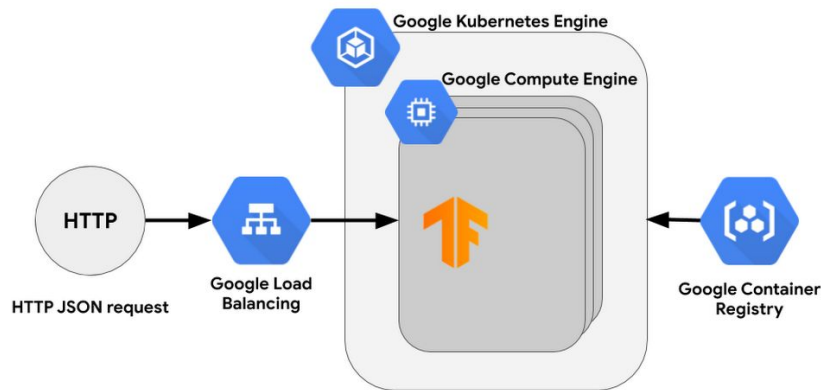bangk!t

# How to **Serve** Your ML Models on **GCP**?

Google Cloud Platform (**GCP**) provides
multiple ways for deploying inference
in the cloud

- [Compute Engine](#)
- [Vertex AI](#)
- [Cloud Functions](#)
- [Cloud Run](#)



**bangkit**

# Scaling ML System with Kubernetes

- In a production setting, you want to be able to scale as the load is increasing on your app
- Kubernetes can help in orchestrating and scaling multiple docker containers



HTTP

HTTP JSON request

Google Load Balancing

Google Kubernetes Engine

Google Compute Engine

Google Container Registry
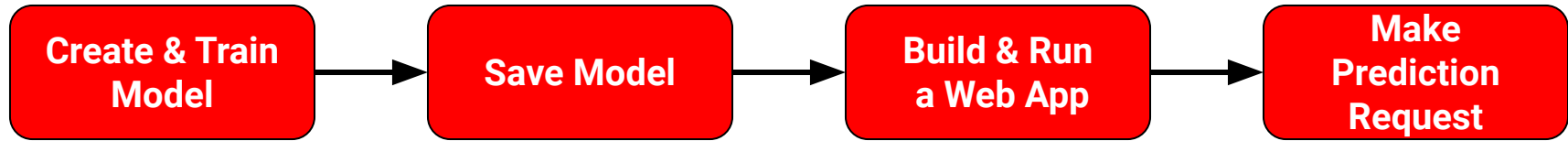
bangkit

# ML Deployment
# Flask

bangkit

# What is **Flask**?

Flask is a minimalist web application framework written in Python

- Easy to use
- Lightweight



bangkit

# Deploy Model with Flask

```
Create & Train Model  →  Save Model  →  Build & Run a Web App  →  Make Prediction Request
```

bangk!t

# How to Build & Run a **Web App**?

**A Simple Flask Web App** (main.py)

```python
from flask import Flask

app = Flask(__name__)

@app.route("/")

def hello_world():

    return "Hello, World!"
```

→ Define the web application

→ Define the route

### Run the Web App

```
export FLASK_APP=main.py

flask run
```

bangkit

# How to Make Prediction Request?

This example demonstrates how to make a prediction route using Flask.

```python
model = joblib.load("iris_model.joblib")
@app.route("/predict", methods=["POST"])
def predict():
    request_json = request.json
    prediction = model.predict(request_json.get("data"))
    prediction_string = [str(d) for d in prediction]
    response_json = {
        "data": request_json.get("data"),
        "prediction": list(prediction_string)
    }
    return json.dumps(response_json)
```

bangk!t

# How to Make Prediction Request?

This example demonstrates how to make a request prediction

```python
import requests
import json


json_data = json.dumps({"data": [[4.9, 3.0, 1.4, 0.2]]})
endpoint = "http://localhost:5000/predict"
headers = {"content-type": "application/json"}
response = requests.post(endpoint, data=json_data, headers=headers)
print(response.json())
```

bangk!t

# Data **Pipelines**

# What is Data Pipeline?



Extract → Transform → Load

bangkit

# TensorFlow **Datasets**

- TensorFlow datasets (TFDS) is a tool provided by TensorFlow to build the ETL process with a consistent API
- TFDS also contains a collection of public research datasets of various types such as audio, text, image, video, etc

**Load**
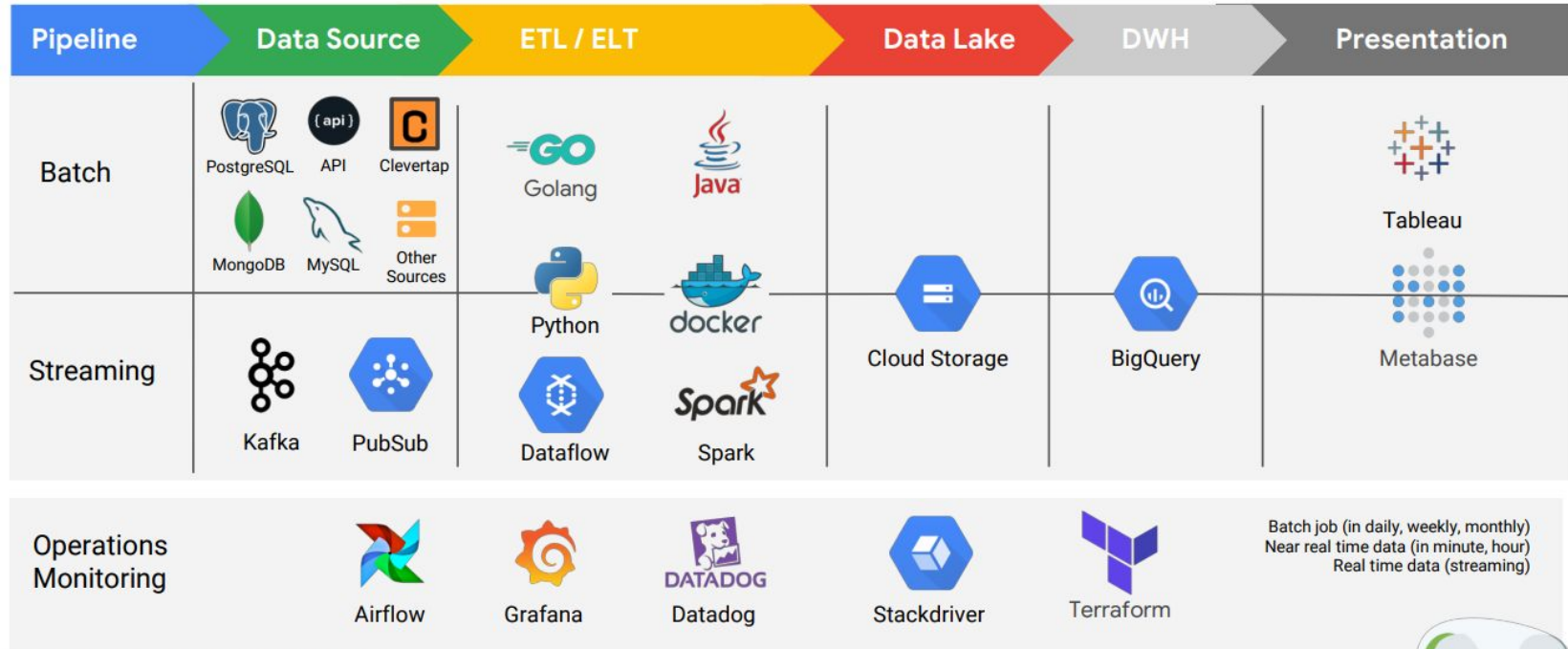
```
For data in dataset.take(10):

…
```

bangk!t

# Gojek's Data Warehouse Architecture



Gojek data architecture as of Q1 2019
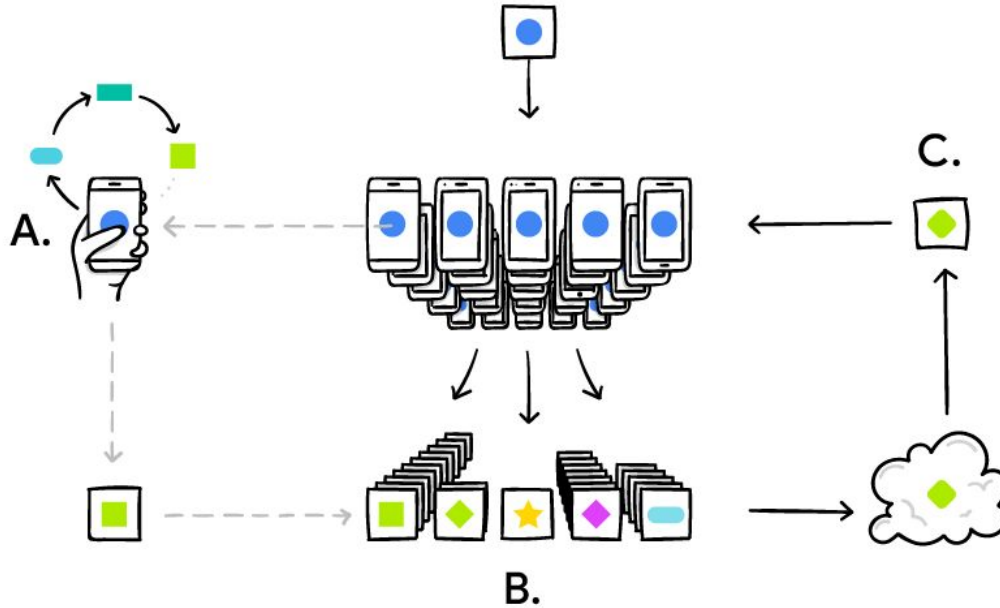
# Federated Learning

# 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



**Federated Learning**

Building better products with on-device data and privacy by default

An online comic from Google AI

bangk!t

# How **Federated Learning** Works?

Image source

# TensorFlow Federated

TensorFlow Federated (TFF) is an open-source framework for machine learning and other computations on decentralized data

- Federated Learning (FL) API
- Federated Core (FC) API

```python
import tensorflow_federated as tff
import nest_asyncio
nest_asyncio.apply()


@tff.federated_computation
def hello_world():
  return "Hello, World!"


hello_world()
```

bangk!t

# Deployment Option **Summary**

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

bangkit

# Sharing Session

bangk!t

# Demo Link

Demo deployment use TF Serving:
https://github.com/dicodingacademy/demo-ilt-ml-bangkit/tree/main/ILT-5/deploy-tf-serving

Demo deployment use Flask:
https://github.com/dicodingacademy/demo-ilt-ml-bangkit/tree/main/ILT-5/deploy-flask

Demo deployment use TFDS:
https://colab.research.google.com/drive/1PnOrYjooTPAa9N4bmejpAxdlLjqZJsGG?usp=sharing

bangkit

# Quiz

# Discussions

bangkit

# Thank You

bangk!t