# Impacts of climate change on food Security
# Group: 17

**Group Members Names:**

1. Mohamed Fadlelseed
2. Fulchany Arvelin Nanitelamio
3. Bacoro dit Elhadji lansar
4. Madalena Makiesse
5. Zeinab Mohmmed
6. Israa Abdulrahman Mohammed kheir

## Deployment

### 1. Overview

The deployment phase is crucial for making the machine learning model for predicting crop yields and assessing the impacts of climate change and undernourishment available in a real-world environment. Key steps include model serialization, serving the model, integrating it with an API for user access, implementing security measures, and setting up monitoring for ongoing performance evaluation.

### 2. Model Serialization

The trained model is serialized to ensure it can be stored and loaded efficiently during deployment. For this project, we utilized the **Pickle** format due to its simplicity and compatibility with Python, which allows for seamless integration with our development environment.

**Serialization Process:**
1. **Model Training Completion**: After the model has been trained and validated, it is prepared for serialization. This involves ensuring that the model's parameters and structure are finalized.
2. **Pickle Serialization**: The model is serialized using the `pickle.dump()` method, which converts the model object into a byte stream that can be saved to disk.

**Considerations for Efficient Storage:**
- **Compression**: To reduce the file size of the serialized model, we apply **gzip compression**. This is achieved by wrapping the Pickle serialization process with the `gzip` module, which significantly decreases the storage requirements while maintaining the integrity of the model data

- **Version Control**: Keeping track of model versions is critical for maintaining an organized deployment process. We implement a versioning system that labels each model file with a version number . This approach allows for easy updates and rollbacks if a newer model version underperforms or introduces issues. Additionally, a metadata file is maintained alongside the serialized model, documenting the training parameters, data used, and performance metrics.

- **Testing Serialized Model**: After serialization, we implement a testing step to ensure that the model can be deserialized successfully and produce accurate predictions. This involves loading the model back into the environment and running a set of validation data through it to verify its integrity and functionality.

- **Compatibility Considerations**: When choosing Pickle, it's important to consider compatibility across different Python versions and environments. To mitigate potential issues, we document the Python version and any library dependencies used during model training, ensuring that the deployment environment mirrors the development setup as closely as possible.

This serialization strategy not only facilitates efficient storage and retrieval of the model but also enhances the overall robustness of the deployment process.

# 3. Model Serving

The serialized model is served using a cloud platform to leverage scalability and reliability, which are crucial for handling variable loads and ensuring consistent performance. For this project, we chose **Docker** to containerize the application, allowing us to deploy the model in a consistent environment.

**Serving Process:**
1. **Docker Container Setup**: We created a Docker image that includes the necessary dependencies, the serialized model, and the code required to serve predictions. This ensures that the model runs in an isolated environment, eliminating issues related to compatibility with different systems.

2. **Container Deployment**: The Docker container is deployed on a local server. This setup allows for easy scaling, as multiple instances of the container can be run concurrently based on demand.

3. **Handling Predictions**: Upon receiving a request, the Docker container processes the input data, which includes various climate variables (temperature, precipitation, soil type) necessary for making predictions about crop yields. The input data is typically in JSON format, making it easy to parse and validate.

4. **Prediction Execution**: The containerized application loads the serialized model into memory and generates predictions based on the incoming data. This architecture ensures quick response times, as the model is readily available within the container.

5. **Returning Results**: After generating predictions, the application returns the results in JSON format, including predicted crop yields along with relevant metadata (confidence intervals). The response is sent back to the client application, depending on the architecture.

**Benefits of Using Docker:**

- **Consistency**: Docker ensures that the application runs the same way across different environments, reducing "works on my machine" issues.
- **Scalability**: Docker containers can be easily replicated and orchestrated using tools like Kubernetes or AWS ECS, allowing for horizontal scaling based on traffic.
- **Isolation**: Each container operates in its own environment, minimizing conflicts

between dependencies and making it easier to manage updates and rollbacks.

## 4. API Integration

The machine learning model is integrated into a **RESTful API** to provide seamless access for users and applications, facilitating the use of predictions in various contexts. This integration is essential for enabling real-time decision-making based on the model's output.

**Key Details:**
**- API Endpoints**:

- A primary endpoint, `/predict`, is created for handling prediction requests regarding crop yields. This endpoint accepts incoming HTTP POST requests with the necessary data for prediction.
- Additional endpoints may be implemented for functionalities such as model health checks (`/health`), version information (`/version`),

or retrieving historical predictions (`/history`).

**- Input Formats**:

- The expected input format is **JSON**, which allows for easy data transmission and parsing. The JSON object should include the following fields:
  - `temperature`: The average temperature expected during the growing season.
  - `precipitation`: The expected rainfall amounts.
  - `soil_type`: A categorical variable indicating the type of soil (e.g., sandy, loamy, clay).
  - `historical_yields`: An array of previous crop yields for the same crop type, used to inform predictions.

**Response Formats**:

- The output format is also **JSON**, containing the predicted yield value and associated metadata. This response provides clarity on the prediction and its reliability. The response object includes:
  - `predicted_yield`: The estimated crop yield based on the input variables.
  - `confidence_interval`: A range indicating the uncertainty of the prediction ( eg. lower and upper bounds).
  - `timestamp`: The time at which the prediction was made, useful for logging and auditing.

**- Error Handling**:

- The API is designed to handle errors gracefully, returning meaningful HTTP status codes and messages for invalid requests. For example:
  - **400 Bad Request**: When the input data does not meet the expected format or is missing required fields.
  - **500 Internal Server Error**: When there is an issue with processing the request, typically indicating a problem with the model or server.

## 5. Security Considerations

To protect sensitive data and ensure safe API access, several security measures are implemented during the deployment phase of the machine learning model. These measures are crucial for maintaining the integrity of the system and safeguarding user information.

**- Authentication**:
- **JWT (JSON Web Tokens)** are utilized to verify user identity before granting access to the API. This method allows users to log in and receive a token that must be included in the header of subsequent requests.
- Tokens contain encoded user information and are signed to prevent tampering. They can also include expiration times to enhance security, requiring users to re-authenticate periodically.

**- Authorization**:

- **Role-Based Access Control (RBAC)** is implemented to ensure that only authorized users can perform specific actions, such as making predictions or accessing sensitive data. Each user is assigned a role (e.g., admin, user, guest) that defines their permissions within the API.
- This approach minimizes the risk of unauthorized access by enforcing strict access policies based on user roles. For example, only admin users may have the ability to manage model versions or access detailed performance metrics.

By implementing these security measures, the deployment phase enhances the overall integrity of the machine learning model and its API, providing a secure environment for users to access predictions related to crop yields and climate impacts.

## 6. Monitoring and Logging

Monitoring the deployed model is essential for maintaining performance and reliability. This includes:
- **Metrics Tracking**: Key performance metrics such as prediction accuracy, request latency, and system resource usage are monitored using tools like AWS CloudWatch.
- **Alerting Mechanisms**: Alerts are set up for significant deviations in performance metrics, ensuring proactive response to potential issues.
- **Logging**: API request and response logs are maintained for auditing and debugging, which helps in assessing model performance and understanding user interactions.