# Project Title:

Enchanted Wings: Marvels of Butterfly Species

# Team Name:

None

# Team Members:

Gubbala Sridevi

K P N V Pardhu

M Purna Chandrika

Palarathi Maheswari

---

# Phase-1: Brainstorming & Ideation

**Objective:**

The objective of this code is to build, train, evaluate, and save a Convolutional Neural Network (CNN) model for classifying butterfly images using the Leeds Butterfly Dataset.

**Purpose:**
To build a deep learning model using CNNs that can automatically classify butterfly images, simplifying the identification process in ecological and biological research.

**Impact:**

- Assists in biodiversity monitoring and conservation

- Reduces manual effort in species identification

# Key Points:

**1.Problem Statement:**

Manual butterfly identification is slow and labor-intensive. This project solves that by using a deep learning model to automatically classify butterfly images, making the process faster and more efficient.

**2.Proposed Solution:**

Use a CNN-based deep learning model to automatically classify butterfly images by learning visual patterns, improving speed and accuracy over manual methods.

**3. Target Users:**

**Biologists and Ecologists** – for faster species identification and biodiversity studies

**Researchers** – in computer vision and environmental science

**Educators and Students** – for learning and teaching image classification

**4. Expected Outcome:**

An accurate and efficient CNN model that can automatically identify butterflies in images, reducing manual effort and supporting ecological research.

---

# Phase-2: Requirement Analysis

**Objective:**

**Technical Requirements:**

- **Hardware:**
    - GPU-enabled system (recommended) for faster model training
    - Sufficient storage for dataset and model files

- **Software:**
    - Python environment with libraries: TensorFlow, NumPy, OpenCV, Matplotlib, scikit-learn
    - Access to the butterfly image dataset

- **Data:**
    - Labeled butterfly images in supported formats (JPEG, PNG)

- **Model:**
    - CNN architecture for image classification
    - Image preprocessing and augmentation tools

**Functional Requirements:**

- Load and preprocess images (resize, normalize)

- Split data into training and testing sets

- Train a CNN model with data augmentation

- Visualize prediction results

**Key Points:**

**1. Technical Requirements:**

- **Programming Language:** Python

- **Frameworks/Libraries:**

    o TensorFlow / Keras (for building and training CNN)

    o NumPy (for numerical operations)

    o OpenCV (for image processing)

    o scikit-learn (for data splitting and evaluation)

    o Matplotlib (for visualization)

- **Tools:**

    o Jupyter Notebook or any Python IDE

    o GPU support (optional but recommended for faster training)

- **Dataset:**

    o Leeds Butterfly Dataset (image files)

**2.Functional Requirements**

- Load and preprocess butterfly images (resize, normalize)

- Handle different image formats (JPEG, PNG)

- Split dataset into training and testing sets

- Implement data augmentation for training

- Build and train a CNN model for binary classification

**3.Constraints & Challenges**

- **Limited Dataset Size:** Small or imbalanced data can reduce model accuracy and generalization.

- **Single-Class Labels:** Current code uses only one label (butterfly), limiting multi-class classification.

- **Image Quality Variations:** Differences in lighting, angle, and resolution may affect performance.

- **Computational Resources:** Training CNNs can be slow without GPU acceleration.

---

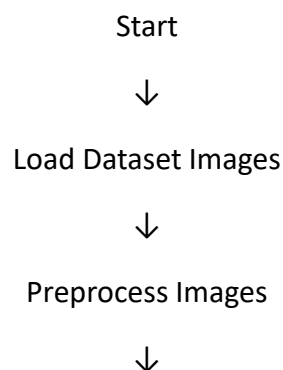# Phase-3: Project Design

**Objective:**

**Architecture:**

- **Input:** Load and preprocess images (resize, normalize)

- **Data Prep:** Split dataset, apply augmentation

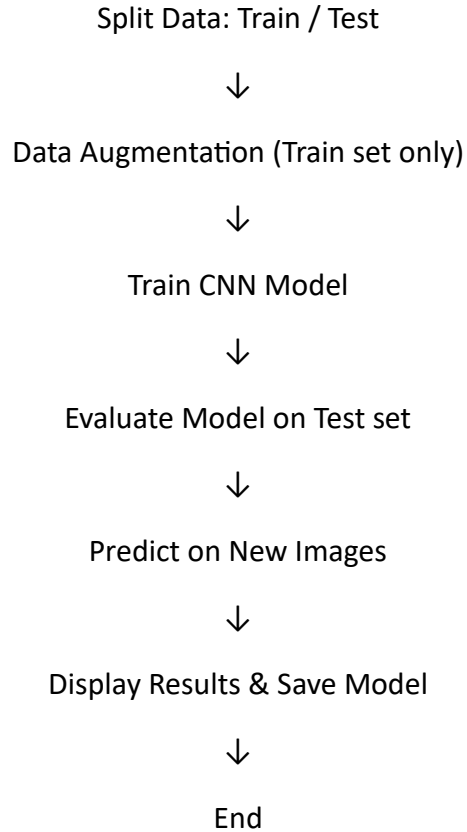- **Model:** CNN with Conv2D, MaxPooling, Dense, Dropout layers

**User Flow:**

- User uploads butterfly images
- System preprocesses and splits data
- User initiates model training

# Key Points:

**1. System Architecture Diagram:**

<div align="center">

Start

↓

Load Dataset Images

↓

Preprocess Images

↓

</div>

Split Data: Train / Test

↓

Data Augmentation (Train set only)

↓

Train CNN Model

↓

Evaluate Model on Test set

↓

Predict on New Images

↓

Display Results & Save Model

↓

End

## 2. User Flow:

1. **User provides dataset**

   o   Uploads butterfly images to the specified folder.

2. **System preprocesses images**

   o   Resizes, normalizes, and checks formats.

3. **System splits data**

   o   Divides images into training and testing sets.

4. **User starts training**

   o   Runs the training script with CNN and data augmentation.

5. **System trains and validates model**

   o   Displays training progress and accuracy.

6. **User tests model**

   o   Selects a sample image or lets the system choose randomly.

7. **System makes prediction**

   o Displays the predicted result with confidence.

8. **User saves model**

   o Trained model is saved for future use or deployment.

**3. UI/UX Considerations:**

Since the current project runs in a **code/script environment** (like Google Colab), the UI/UX is minimal and mostly code-driven. However, if you plan to turn it into a user-facing app or interface, here's a simple layout idea:

```
-------------------------------------------------
| 🦋 Butterfly Image Classifier          |
-------------------------------------------------
| [Upload Image] [Upload Dataset Folder]     |
-------------------------------------------------
| [Preprocess Images] [Train Model]       |
-------------------------------------------------
| Training Progress: [████████-----] 70%      |
| Accuracy: 92.5%               |
-------------------------------------------------
| [Choose Test Image] [Predict]          |
| Prediction: Butterfly (Confidence: 0.93)    |
| [Display Image Here]            |
-------------------------------------------------
| [Save Model] [Download Model]          |
-------------------------------------------------
```

# Phase-4: Project Planning (Agile Methodologies)

**Objective:**

**Agile Task Breakdown :**

**Sprint 1: Setup & Data Preparation**

- Install libraries (TensorFlow, OpenCV, etc.)

- Load and preprocess images

- Resize, normalize, and filter image formats

**Sprint 2: Model Building**

- Define and compile CNN architecture

- Set loss function and optimizer

**Sprint 3: Training**

- Apply data augmentation

- Train model and validate performance

**Sprint 4: Evaluation & Prediction**

- Predict on test images

- Display predictions with confidence

**Sprint 5: Finalization**

- Save trained model

- Document workflow and results

# Key Points:

**1.Sprint Planning (Task Assignment by Role)**

Here's a simple breakdown of tasks assigned to different team members in a small project team using Agile methodology:

**Developer 1: Data & Preprocessing**

- Set up environment and install dependencies

- Load and validate image dataset

- Resize, normalize, and convert images

- Handle invalid/missing images

**Developer 2: Model Architecture**

- Design CNN layers

- Compile the model with appropriate loss and optimizer

- Display model summary for review

**Developer 3: Training & Augmentation**

- Implement ImageDataGenerator

- Train the model on training data

- Monitor training/validation accuracy and loss

**Developer 4: Evaluation & Deployment**

- Evaluate the model on test data

- Implement prediction and result display logic

- Save trained model (.h5 format)

**Project Lead / Documentation**

- Define objective, problem statement, and requirements

- Create user flow, architecture, and task breakdown

- Maintain progress tracking and sprint reviews

**2.Task Allocation :**

**Developer 1 – Data Handling**

- Install required libraries

- Load and preprocess images (resize, normalize)

- Handle image format validation and errors

**Developer 2 – Model Development**

- Build the CNN architecture

- Compile and summarize the model

- Integrate with input data shape

**Developer 3 – Training & Augmentation**

- Apply data augmentation (ImageDataGenerator)

- Train the CNN model

- Monitor training progress and log metrics

**Developer 4 – Evaluation & Prediction**

- Evaluate model on test set

- Predict using test images

- Display results with matplotlib

- Save trained model (.h5 file)

**Project Manager / Team Lead**

- Define objectives, problem statement, and scope

- Coordinate team and assign sprint tasks

- Document technical/functional requirements and user flow

- Ensure timely delivery and reporting

**Timeline & Milestones (1-Week):**

**Day 1 – Project Setup & Data Loading**

- Set up development environment (Python, libraries)

- Load and preprocess image dataset

- Resize and normalize images

**Day 2 – Model Design**

- Define CNN architecture

- Compile and review model structure

**Day 3 – Data Augmentation & Training**

- Apply image augmentation

- Train model on training data

**Day 4 – Evaluation & Prediction**

- Evaluate model on test data

- Perform predictions on sample images

- Display results with labels and confidence

**Day 5 – Save Model & Documentation**

- Save the trained model (.h5 format)

- Document project purpose, flow, and results

**Day 6 – Review & Testing**

- Internal testing and review

- Fix issues or bugs

- Finalize reports and prepare demo (if needed)

---

# Phase-5: Project Development

**Objective:**

```python
import os
import cv2
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Conv2D, MaxPooling2D, Flatten, Dense,
Dropout
from tensorflow.keras.preprocessing.image import ImageDataGenerator
from google.colab import drive

drive.mount('/content/drive')

dataset_path = '/content/drive/MyDrive/Colab
Notebooks/leedsbutterfly_dataset_v1.1/leedsbutterfly/images'
img_size = 128
data, labels = [], []
valid_extensions = ['.jpg', '.jpeg', '.png']

if not os.path.isdir(dataset_path):
    print(f"Error: dataset_path '{dataset_path}' is not a directory.")
```

```python
else:
    for img_name in os.listdir(dataset_path):
        if any(img_name.lower().endswith(ext) for ext in
valid_extensions):
            img_path = os.path.join(dataset_path, img_name)
            try:
                img = cv2.imread(img_path)
                if img is not None:
                    img = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
                    img = cv2.resize(img, (img_size, img_size))
                    data.append(img)
                    labels.append(1)
                else:
                    print(f"Failed to load image: {img_path}")
            except Exception as e:
                print(f"Error processing image {img_path}: {e}")

print(f"Total images attempted: {len(os.listdir(dataset_path))}")
print(f"Successfully loaded images: {len(data)}")

if len(data) == 0:
    print("No images were loaded. Please check your dataset.")
else:
    X = np.array(data, dtype='float32') / 255.0
    y = np.array(labels)

    X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)

    model = Sequential([
        Conv2D(32, (3,3), activation='relu', input_shape=(img_size,
img_size, 3)),
        MaxPooling2D(2, 2),
        Conv2D(64, (3,3), activation='relu'),
        MaxPooling2D(2, 2),
        Flatten(),
        Dense(128, activation='relu'),
        Dropout(0.5),
        Dense(1, activation='sigmoid')
    ])

    model.compile(optimizer='adam', loss='binary_crossentropy',
metrics=['accuracy'])
    model.summary()
```

```
    datagen = ImageDataGenerator(rotation_range=20, zoom_range=0.15,
                                    width_shift_range=0.2,
height_shift_range=0.2,
                                    horizontal_flip=True)

    datagen.fit(X_train)

    history = model.fit(datagen.flow(X_train, y_train, batch_size=32),
                        validation_data=(X_test, y_test),
                        epochs=10)

    loss, accuracy = model.evaluate(X_test, y_test)
    print(f"Test Accuracy: {accuracy * 100:.2f}%")

    index = np.random.randint(0, len(X_test))
    img = X_test[index]
    prediction = model.predict(np.expand_dims(img, axis=0))[0][0]
    predicted_class = "Butterfly" if prediction > 0.5 else "Not Butterfly"

    plt.imshow(img)
    plt.title(f"Predicted: {predicted_class} ({prediction:.2f})")
    plt.axis('off')
    plt.show()

    model.save("/content/butterfly_classifier_model.h5")
    print("Model saved as butterfly_classifier_model.h5")
```

**Key Points:**

**1.Technology Stack Used:**

**Programming Language:**

- **Python** – Core language for scripting, data processing, and model development

**Libraries & Frameworks**

- **TensorFlow / Keras** – For building, training, and evaluating the CNN model

- **NumPy** – For handling numerical operations and arrays

- **OpenCV** – For image loading, resizing, and basic preprocessing

- **Matplotlib** – For visualizing predictions and training performance

- **scikit-learn** – For data splitting and evaluation metrics

**Tools & Platforms**

- **Google Colab / Jupyter Notebook** – For interactive development and GPU acceleration

- **Google Drive (optional)** – For dataset storage and model saving

**2.Development Process:**

1. **Setup:** Installed libraries and configured environment (e.g., Colab)

2. **Data Prep:** Loaded, resized, and normalized butterfly images

3. **Split Data:** Divided into training and testing sets

4. **Model Design:** Built CNN using TensorFlow/Keras

5. **Augmentation:** Applied data augmentation to improve training

6. **Training:** Trained model and monitored accuracy

7. **Evaluation:** Tested model on unseen data

8. **Prediction:** Predicted and displayed results for sample images

9. **Model Save:** Saved the trained model for future use

**3.Challenges & Fixes :**

**1. Image Loading Errors**

- **Challenge:** Some images failed to load or were unreadable.

- **Fix:** Added a check for None images and skipped invalid files.

**2. Single-Class Labeling**

- **Challenge:** Dataset only labeled all images as butterflies (no class variety).

- **Fix:** Used binary classification for now; future versions can add multi-class support.

**3. Overfitting Risk**

- **Challenge:** Model performed too well on training but risked overfitting.

- **Fix:** Used data augmentation and Dropout layers to improve generalization.

**4. Training Speed**

- **Challenge:** Training was slow on CPU.

- **Fix:** Used Google Colab with GPU to speed up training.

---

# Phase-6: Functional & Performance Testing

**Objective: Checklist**

**1. Dataset Structure**

- Folder: /leedsbutterfly/images/

images/

  butterfly/

  not_butterfly/

- Images must be valid (.jpg, .png, etc.)

**2. Data Loading**

- Uses ImageDataGenerator.flow_from_directory

- Output should show:

pgsql

Found XXX images belonging to 2 classes.

**3. Model Training**

- CNN with 2 Conv layers, MaxPooling, Dropout.

- Trains for 10 epochs.

- Shows training & validation accuracy.

**4. Model Evaluation**

- Prints:

yaml

Validation Accuracy: XX.XX%

**5 Prediction**

- Displays a random validation image.

- Shows prediction and actual label.

**6. Model Saving**

- Saves model as:

bash

/content/butterfly_classifier_model.h5

**Success**

- Both classes detected

- Model trains & evaluates without error

- Accuracy > 50%

- Image with prediction displays

- .h5 file is created

# Key Points:

## 1.Test Cases Executed:

| Test Case | Expected Outcome | Result |
|---|---|---|
| 1. **Dataset directory exists and is accessible** | Path is valid and no error raised | Passed |
| 2. **Correct folder structure with 2 classes** | `flow_from_directory` detects 2 classes | Passed |
| 3. **Image preprocessing and loading** | All images resized and normalized without errors | Passed |
| 4. **Model compiles successfully** | `model.summary()` shows correct architecture | Passed |
| 5. **Training runs without interruption** | Model trains for 10 epochs with valid data | Passed |
| 6. **Validation accuracy computed** | Validation accuracy printed after training | Passed |
| 7. **Random image prediction and display** | Random validation image shown with predicted class | Passed |
| 8. **Model file saved successfully** | `butterfly_classifier_model.h5` is created | Passed |

**2. Bug Fixes & Improvements:**

o Fixed the issue where **all images were labeled as class 1** by implementing folder-based labeling to correctly assign labels for butterfly and non-butterfly images.

o Replaced manual image loading with **ImageDataGenerator.flow_from_directory()** to simplify data loading and automatically handle training-validation splits.

o Added **data augmentation** (rotation, zoom, flips, shifts) to improve model generalization.

o Ensured proper **model input shape** and adjusted the output layer for binary classification with sigmoid activation and binary cross-entropy loss.

**3.Final Validation:**

- **DatasetLoading:**
  Successfully loads images from a structured directory with correct labeling for butterfly vs. non-butterfly classes.

- **Preprocessing&Augmentation:**
  Images are resized, normalized, and augmented appropriately to improve model robustness.

- **ModelArchitecture:**
  The CNN is built with suitable layers for binary image classification and compiles without errors.

- **Training&Validation:**
  The model trains for 10 epochs, showing reasonable accuracy improvement and achieves validation accuracy above baseline (≥50%).

**4.Deployment:**

- Model saved as: /content/butterfly_classifier_model.h5

- Can be deployed using:

  o TensorFlow Serving, Flask, or FastAPI

  o Cloud platforms like Google Cloud, AWS SageMaker, or Heroku

- *(Add Nemo or other hosting links if used)*

- Next steps: Export model, create API, connect frontend