```python
# This Python 3 environment comes with many helpful analytics libraries installed
# It is defined by the kaggle/python Docker image: https://github.com/kaggle/docker-python
# For example, here's several helpful packages to load

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

# Input data files are available in the read-only "../input/" directory
# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input direct

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

# You can write up to 20GB to the current directory (/kaggle/working/) that gets preserved as output when you cr
# You can also write temporary files to /kaggle/temp/, but they won't be saved outside of the current session
```

```python
# 1. Install Kaggle package
!pip install -q kaggle

# 2. Upload your API token file (kaggle.json)
from google.colab import files
files.upload()  # choose your kaggle.json file

# 3. Move it to correct location & set permissions
!mkdir -p ~/.kaggle
!cp kaggle.json ~/.kaggle/
!chmod 600 ~/.kaggle/kaggle.json

# 4. Download the dataset
!kaggle datasets download -d nigarmahmoudshafiq/nigar-eeg-alzheimers-dataset-v1

# 5. Unzip it (if zipped)
!unzip nigar-eeg-alzheimers-dataset-v1.zip
```

Choose Files | kaggle.json
**kaggle.json**(application/json) - 68 bytes, last modified: 11/4/2025 - 100% done
Saving kaggle.json to kaggle.json
Dataset URL: https://www.kaggle.com/datasets/nigarmahmoudshafiq/nigar-eeg-alzheimers-dataset-v1
License(s): Attribution 4.0 International (CC BY 4.0)
Downloading nigar-eeg-alzheimers-dataset-v1.zip to /content
 96% 1.67G/1.74G [00:12<00:01, 48.1MB/s]
100% 1.74G/1.74G [00:12<00:00, 149MB/s]
Archive:  nigar-eeg-alzheimers-dataset-v1.zip
  inflating: Final dataset for the published paper/1-Healthy/CN1.csv
  inflating: Final dataset for the published paper/1-Healthy/CN10.csv
  inflating: Final dataset for the published paper/1-Healthy/CN11.csv
  inflating: Final dataset for the published paper/1-Healthy/CN12.csv
  inflating: Final dataset for the published paper/1-Healthy/CN13.csv
  inflating: Final dataset for the published paper/1-Healthy/CN14.csv
  inflating: Final dataset for the published paper/1-Healthy/CN15.csv
  inflating: Final dataset for the published paper/1-Healthy/CN16.csv
  inflating: Final dataset for the published paper/1-Healthy/CN17.csv
  inflating: Final dataset for the published paper/1-Healthy/CN18.csv
  inflating: Final dataset for the published paper/1-Healthy/CN19.csv
  inflating: Final dataset for the published paper/1-Healthy/CN2.csv
  inflating: Final dataset for the published paper/1-Healthy/CN20.csv
  inflating: Final dataset for the published paper/1-Healthy/CN21.csv
  inflating: Final dataset for the published paper/1-Healthy/CN22.csv
  inflating: Final dataset for the published paper/1-Healthy/CN23.csv
  inflating: Final dataset for the published paper/1-Healthy/CN3.csv
  inflating: Final dataset for the published paper/1-Healthy/CN4.csv
  inflating: Final dataset for the published paper/1-Healthy/CN5.csv
  inflating: Final dataset for the published paper/1-Healthy/CN6.csv
  inflating: Final dataset for the published paper/1-Healthy/CN7.csv
  inflating: Final dataset for the published paper/1-Healthy/CN8.csv
  inflating: Final dataset for the published paper/1-Healthy/CN9.csv
  inflating: Final dataset for the published paper/2-Mild/Mild1.csv
  inflating: Final dataset for the published paper/2-Mild/Mild2.csv
  inflating: Final dataset for the published paper/2-Mild/Mild3.csv
  inflating: Final dataset for the published paper/2-Mild/Mild4.csv
  inflating: Final dataset for the published paper/2-Mild/Mild5.csv
  inflating: Final dataset for the published paper/2-Mild/Mild6.csv
  inflating: Final dataset for the published paper/2-Mild/Mild7.csv
  inflating: Final dataset for the published paper/2-Mild/Mild8.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod1.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod10.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod11.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod12.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod2.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod3.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod4.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod5.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod6.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod7.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod8.csv
  inflating: Final dataset for the published paper/3-Moderate/Mod9.csv
  inflating: Final dataset for the published paper/4-Sever/Sv1.csv
  inflating: Final dataset for the published paper/4-Sever/Sv10.csv
  inflating: Final dataset for the published paper/4-Sever/Sv2.csv
  inflating: Final dataset for the published paper/4-Sever/Sv3.csv
  inflating: Final dataset for the published paper/4-Sever/Sv4.csv
  inflating: Final dataset for the published paper/4-Sever/Sv5.csv
  inflating: Final dataset for the published paper/4-Sever/Sv6.csv
  inflating: Final dataset for the published paper/4-Sever/Sv7.csv
  inflating: Final dataset for the published paper/4-Sever/Sv8.csv
  inflating: Final dataset for the published paper/4-Sever/Sv9.csv
  inflating: Final dataset for the published paper/HMMS.csv

```python
import os
import pandas as pd
import numpy as np

# Root directory where the folders are located
root = "/content/Initial_Dataset"   # change if needed

groups = ["Healthy", "Mild", "Moderate", "Severe"]

for grp in groups:
    folder_path = os.path.join(root, grp)

    # Add a check in case the folder doesn't exist
    if not os.path.exists(folder_path):
        print(f"Skipping → {folder_path} (Not found)")
        continue

    for file in os.listdir(folder_path):
        if file.endswith(".csv"):
            file_path = os.path.join(folder_path, file)
            print(f"Processing → {file_path}")
```

```
            df = pd.read_csv(file_path)

            # --- FIX: Use column names for robust assignment ---

            # 1. Get the names of the columns at those positions
            col_names = df.columns[[0, 1, 18]]
            col_fp1, col_fp2, col_pz = col_names

            # 2. Extract the data using those names and cast to float
            selected = df[col_names].astype(float)

            # 3. Rename columns in the new 'selected' DataFrame for clarity
            selected.columns = ["Fp1", "Fp2", "Pz"]

            # 4. Remove 5% of values randomly
            num_values = selected.size
            num_missing = int(num_values * 0.05)

            rows = np.random.randint(0, selected.shape[0], num_missing)
            cols = np.random.randint(0, selected.shape[1], num_missing)

            for r, c in zip(rows, cols):
                selected.iat[r, c] = np.nan # This is fine, 'selected' is already float

            # 5. Write the modified float columns back into the original df
            #    Using the original column names (e.g., df[col_fp1] = ...)
            #    will correctly handle the dtype change without warnings.
            df[col_fp1] = selected["Fp1"]
            df[col_fp2] = selected["Fp2"]
            df[col_pz] = selected["Pz"]

            # 6. Save back (overwrite)
            df.to_csv(file_path, index=False)

    print("\n✅ DONE – Missing values inserted into all subjects.")
```

```
Processing → /content/Initial_Dataset/Healthy/CN6.csv
Processing → /content/Initial_Dataset/Healthy/CN19.csv
Processing → /content/Initial_Dataset/Healthy/CN17.csv
Processing → /content/Initial_Dataset/Healthy/CN11.csv
Processing → /content/Initial_Dataset/Healthy/CN3.csv
Processing → /content/Initial_Dataset/Healthy/CN20.csv
Processing → /content/Initial_Dataset/Healthy/CN8.csv
Processing → /content/Initial_Dataset/Healthy/CN18.csv
Processing → /content/Initial_Dataset/Healthy/CN13.csv
Processing → /content/Initial_Dataset/Healthy/CN12.csv
Processing → /content/Initial_Dataset/Healthy/CN21.csv
Processing → /content/Initial_Dataset/Healthy/CN10.csv
Processing → /content/Initial_Dataset/Healthy/CN14.csv
Processing → /content/Initial_Dataset/Healthy/CN4.csv
Processing → /content/Initial_Dataset/Healthy/CN7.csv
Processing → /content/Initial_Dataset/Healthy/CN9.csv
Processing → /content/Initial_Dataset/Healthy/CN22.csv
Processing → /content/Initial_Dataset/Healthy/CN1.csv
Processing → /content/Initial_Dataset/Healthy/CN2.csv
Processing → /content/Initial_Dataset/Healthy/CN16.csv
Processing → /content/Initial_Dataset/Healthy/CN23.csv
Processing → /content/Initial_Dataset/Healthy/CN15.csv
Processing → /content/Initial_Dataset/Healthy/CN5.csv
Processing → /content/Initial_Dataset/Mild/Mild8.csv
Processing → /content/Initial_Dataset/Mild/Mild7.csv
Processing → /content/Initial_Dataset/Mild/Mild4.csv
Processing → /content/Initial_Dataset/Mild/Mild3.csv
Processing → /content/Initial_Dataset/Mild/Mild5.csv
Processing → /content/Initial_Dataset/Mild/Mild6.csv
Processing → /content/Initial_Dataset/Mild/Mild2.csv
Processing → /content/Initial_Dataset/Mild/Mild1.csv
Processing → /content/Initial_Dataset/Moderate/Mod6.csv
Processing → /content/Initial_Dataset/Moderate/Mod1.csv
Processing → /content/Initial_Dataset/Moderate/Mod5.csv
Processing → /content/Initial_Dataset/Moderate/Mod3.csv
Processing → /content/Initial_Dataset/Moderate/Mod7.csv
Processing → /content/Initial_Dataset/Moderate/Mod11.csv
Processing → /content/Initial_Dataset/Moderate/Mod9.csv
Processing → /content/Initial_Dataset/Moderate/Mod8.csv
Processing → /content/Initial_Dataset/Moderate/Mod4.csv
Processing → /content/Initial_Dataset/Moderate/Mod10.csv
Processing → /content/Initial_Dataset/Moderate/Mod2.csv
Processing → /content/Initial_Dataset/Moderate/Mod12.csv
Processing → /content/Initial_Dataset/Severe/Sv9.csv
Processing → /content/Initial_Dataset/Severe/Sv6.csv
Processing → /content/Initial_Dataset/Severe/Sv8.csv
Processing → /content/Initial_Dataset/Severe/Sv5.csv
Processing → /content/Initial_Dataset/Severe/Sv3.csv
Processing → /content/Initial_Dataset/Severe/Sv7.csv
Processing → /content/Initial_Dataset/Severe/Sv1.csv
```

```
Processing → /content/Initial_Dataset/Severe/Sv4.csv
Processing → /content/Initial_Dataset/Severe/Sv10.csv
Processing → /content/Initial_Dataset/Severe/Sv2.csv

✅ DONE — Missing values inserted into all subjects.
```

```python
# --------------------------------------------------------------------
# STEP 0: IMPORT LIBRARIES
# --------------------------------------------------------------------
import numpy as np
import pandas as pd
import os
import glob
from scipy import signal as sp_signal
from scipy.stats import skew, kurtosis, entropy
from scipy.integrate import simpson
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Sklearn Models (SVM Removed)
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier

# --- Imputation Imports ---
# Install the firefly library
!pip install -q fireflyalgorithm
from fireflyalgorithm import FireflyAlgorithm
from sklearn.impute import KNNImputer
from sklearn.metrics import mean_squared_error

# --- Visualization ---
import matplotlib.pyplot as plt
import seaborn as sns

import warnings
warnings.filterwarnings('ignore')

print("Step 0: All libraries (except SVM) imported successfully.")


# --------------------------------------------------------------------
# STEP 1: DEFINE CONSTANTS AND LOCATE DATA FILES
# --------------------------------------------------------------------
CHANNELS_TO_USE_INDICES = [0, 1, 18] # 1st, 2nd, 19th columns
SAMPLING_RATE_HZ = 250
NUM_CLASSES = 4

# --- Epoching Configuration ---
EPOCH_SEC = 2.0
EPOCH_SAMPLES = int(SAMPLING_RATE_HZ * EPOCH_SEC)
OVERLAP_RATIO = 0.5
STEP_SAMPLES = int(EPOCH_SAMPLES * (1.0 - OVERLAP_RATIO))
MIN_SAMPLES_REQUIRED = EPOCH_SAMPLES

# --- Frequency Bands ---
FREQ_BANDS = {
    'Delta': (1, 4), 'Theta': (4, 8), 'Alpha': (8, 13), 'Beta': (13, 30)
}

# --- File Discovery (Assumes CORRUPTED dataset) ---
BASE_DIR = '/content/Initial_Dataset'
CLASSES = ['Healthy', 'Mild', 'Moderate', 'Severe']
LABEL_MAP = {'Healthy': 0, 'Mild': 1, 'Moderate': 2, 'Severe': 3}

subject_files_by_class = {}
for class_name in CLASSES:
    search_path = os.path.join(BASE_DIR, class_name, '*.csv')
    files = glob.glob(search_path)
    # This dataset has 'CN' for 'Healthy' in some cases
    if class_name == 'Healthy':
        search_path_cn = os.path.join(BASE_DIR, 'CN', '*.csv')
        files.extend(glob.glob(search_path_cn))
    subject_files_by_class[class_name] = list(set(files))
    print(f"Found {len(files)} subjects for class: {class_name}")

if sum(len(v) for v in subject_files_by_class.values()) == 0:
    print("\nCRITICAL ERROR: Could not find any data files.")
    print(f"Please check that BASE_DIR is set to: {BASE_DIR}")
    PROCEDE = False
```

```python
        PROCEDE = False
    else:
        PROCEDE = True
        print("\nStep 1: File discovery complete. (Loading from corrupted source)")



    # -------------------------------------------------------------------------
    # STEP 2: SPLIT SUBJECTS (File Paths)
    # -------------------------------------------------------------------------
    if PROCEDE:
        train_files = []
        test_files = []

        for class_name, files in subject_files_by_class.items():
            if not files:
                continue
            # 85% Train / 15% Test split
            files_train, files_test = train_test_split(
                files, test_size=0.15, random_state=42
            )
            label = LABEL_MAP[class_name]
            train_files.extend([(f, label) for f in files_train])
            test_files.extend([(f, label) for f in files_test])

        print(f"\nTraining subjects: {len(train_files)} (85%)")
        print(f"Testing subjects: {len(test_files)} (15%)")
        print("\nStep 2: Subject-based data splitting complete.")



    # -------------------------------------------------------------------------
    # STEP 3: FEATURE ENGINEERING (NaN-Aware)
    # -------------------------------------------------------------------------
    def calculate_band_powers(data, fs, bands):
        nperseg = len(data)
        freqs, psd = sp_signal.welch(data, fs=fs, nperseg=nperseg)
        total_power = simpson(psd[(freqs >= 1) & (freqs <= 30)], freqs[(freqs >= 1) & (freqs <= 30)])
        powers = {}
        for band, (low, high) in bands.items():
            band_idx = np.logical_and(freqs >= low, freqs <= high)
            abs_power = simpson(psd[band_idx], freqs[band_idx])
            powers[f'Abs_{band}'] = abs_power
            powers[f'Rel_{band}'] = abs_power / (total_power + 1e-10)
        return powers

    def calculate_hjorth(signal):
        activity = np.var(signal) + 1e-10
        diff1 = np.diff(signal)
        var_diff1 = np.var(diff1) + 1e-10
        mobility = np.sqrt(var_diff1 / activity)
        diff2 = np.diff(diff1)
        var_diff2 = np.var(diff2) + 1e-10
        complexity = np.sqrt(var_diff2 / var_diff1) / mobility
        return mobility, complexity

    def process_subject_file_features(file_path, label):
        """
        Processes one *corrupted* subject file.
        If an epoch contains NaN, all its features will be NaN.
        """
        try:
            df = pd.read_csv(file_path)
            if len(df) < MIN_SAMPLES_REQUIRED: return []
            max_col_index = max(CHANNELS_TO_USE_INDICES)
            if len(df.columns) <= max_col_index: return []

            cols_to_use = df.columns[CHANNELS_TO_USE_INDICES]
            channel_data = {name: df[name].values for name in cols_to_use}
            total_samples = len(df)
            subject_epochs = []

            for start in range(0, total_samples - EPOCH_SAMPLES + 1, STEP_SAMPLES):
                end = start + EPOCH_SAMPLES
                epoch_features = {'label': label}
                segments = {}

                all_segments_clean = True

                for channel_name in cols_to_use:
                    segment = channel_data[channel_name][start:end]

                    # --- THIS IS THE NEW, CRITICAL LOGIC ---
                    if np.isnan(segment).any():
                        all_segments_clean = False # Mark this epoch as corrupted
                        # Create NaN entries for all features for this channel
```

```python
                # create NaN entries for all features for this channel
                epoch_features[f'{channel_name}_Alpha_Beta_Ratio'] = np.nan
                epoch_features[f'{channel_name}_mean'] = np.nan
                epoch_features[f'{channel_name}_variance'] = np.nan
                epoch_features[f'{channel_name}_skewness'] = np.nan
                epoch_features[f'{channel_name}_kurtosis'] = np.nan
                epoch_features[f'{channel_name}_shannon_entropy'] = np.nan
                epoch_features[f'{channel_name}_hjorth_mobility'] = np.nan
                epoch_features[f'{channel_name}_hjorth_complexity'] = np.nan
                for band in FREQ_BANDS:
                    epoch_features[f'{channel_name}_Abs_{band}'] = np.nan
                    epoch_features[f'{channel_name}_Rel_{band}'] = np.nan
            else:
                # Segment is clean, store it for processing
                segments[channel_name] = segment

        # Now, process the clean segments (if all were clean)
        if all_segments_clean:
            for channel_name in cols_to_use:
                segment = segments[channel_name]
                band_powers = calculate_band_powers(segment, SAMPLING_RATE_HZ, FREQ_BANDS)
                for key, val in band_powers.items():
                    epoch_features[f'{channel_name}_{key}'] = val
                epoch_features[f'{channel_name}_Alpha_Beta_Ratio'] = band_powers.get('Abs_Alpha', 0) / (band_
                epoch_features[f'{channel_name}_mean'] = np.mean(segment)
                epoch_features[f'{channel_name}_variance'] = np.var(segment)
                epoch_features[f'{channel_name}_skewness'] = skew(segment)
                epoch_features[f'{channel_name}_kurtosis'] = kurtosis(segment)
                epoch_features[f'{channel_name}_shannon_entropy'] = entropy(np.histogram(segment, bins=30, de
                mobility, complexity = calculate_hjorth(segment)
                epoch_features[f'{channel_name}_hjorth_mobility'] = mobility
                epoch_features[f'{channel_name}_hjorth_complexity'] = complexity

            # Correlation
            c_names = list(cols_to_use)
            epoch_features['corr_0_1'] = np.nan_to_num(np.corrcoef(segments[c_names[0]], segments[c_names[1]]
            epoch_features['corr_0_2'] = np.nan_to_num(np.corrcoef(segments[c_names[0]], segments[c_names[2]]
            epoch_features['corr_1_2'] = np.nan_to_num(np.corrcoef(segments[c_names[1]], segments[c_names[2]]

        else: # If *any* segment was NaN, correlation is also NaN
            epoch_features['corr_0_1'] = np.nan
            epoch_features['corr_0_2'] = np.nan
            epoch_features['corr_1_2'] = np.nan

        subject_epochs.append(epoch_features)
    return subject_epochs
    except Exception as e:
        return []

def create_dataset_from_files_features(file_list):
    """Creates the feature DataFrame (will contain NaNs)."""
    features_list = []
    for i, (file_path, label) in enumerate(file_list):
        if (i + 1) % 5 == 0 or i == 0 or i == len(file_list) - 1:
            print(f"  [Feature] Processing subject {i+1}/{len(file_list)}")
        features_from_subject = process_subject_file_features(file_path, label)
        if features_from_subject:
            features_list.extend(features_from_subject)
    return pd.DataFrame(features_list)

print("\nStep 3: NaN-aware feature engineering pipeline defined.")


# ----------------------------------------------------------------------
# STEP 4: CREATE (NaN-filled) DATASETS
# ----------------------------------------------------------------------
if PROCEDE:
    print("\n--- Creating Feature Datasets (will contain NaNs) ---")
    df_train = create_dataset_from_files_features(train_files)
    df_test = create_dataset_from_files_features(test_files)

    if df_train.empty or df_test.empty:
        print("\nCRITICAL ERROR: No data was loaded. Cannot proceed.")
        PROCEDE = False
    else:
        # X_train and X_test will have NaN values
        X_train = df_train.drop(columns=['label'])
        y_train = df_train['label']

        X_test = df_test.drop(columns=['label'])
        y_test = df_test['label']

        # Align columns
```

```
        X_test = X_test.reindex(columns=X_train.columns, fill_value=0)

        print(f"Training data shape (pre-imputation): {X_train.shape}")
        print(f"Test data shape (pre-imputation): {X_test.shape}")

        # Check how many NaNs we created
        train_nans = X_train.isna().sum().sum()
        test_nans = X_test.isna().sum().sum()
        print(f"Found {train_nans} NaN feature values in training set.")
        print(f"Found {test_nans} NaN feature values in test set.")
```

```
Step 0: All libraries (except SVM) imported successfully.
Found 23 subjects for class: Healthy
Found 8 subjects for class: Mild
Found 12 subjects for class: Moderate
Found 10 subjects for class: Severe

Step 1: File discovery complete. (Loading from corrupted source)

Training subjects: 43 (85%)
Testing subjects: 10 (15%)

Step 2: Subject-based data splitting complete.

Step 3: NaN-aware feature engineering pipeline defined.

--- Creating Feature Datasets (will contain NaNs) ---
  [Feature] Processing subject 1/43
  [Feature] Processing subject 5/43
  [Feature] Processing subject 10/43
  [Feature] Processing subject 15/43
  [Feature] Processing subject 20/43
  [Feature] Processing subject 25/43
  [Feature] Processing subject 30/43
  [Feature] Processing subject 35/43
  [Feature] Processing subject 40/43
  [Feature] Processing subject 43/43
  [Feature] Processing subject 1/10
  [Feature] Processing subject 5/10
  [Feature] Processing subject 10/10
Training data shape (pre-imputation): (47684, 115)
Test data shape (pre-imputation): (13224, 115)
Found 5483660 NaN feature values in training set.
Found 1309176 NaN feature values in test set.
```

```python
# ------------------------------------------------------------------------
# STEP 5: IMPUTATION WITH FIREFLY + KNN
# ------------------------------------------------------------------------
if PROCEDE:
    print("\n--- Step 5: Running Firefly + KNN Imputation ---")

    # We must find the best 'k' using *only* the training data

    # Fitness function for FA
    # We will run it on the *training set*.
    def fitness(k_array):
        k = int(np.clip(np.round(k_array[0]), 1, 15))
        imputer = KNNImputer(n_neighbors=k)

        try:
            imputed = imputer.fit_transform(X_train)
            # Return the sum of variances of the imputed data
            # This is a proxy for "stability"
            return np.nanvar(imputed)
        except Exception:
            return 1e10 # Penalize failure

    # --- FIX IS HERE ---
    # 'n_iterations' is removed from the constructor.
    fa = FireflyAlgorithm(pop_size=20)

    print("\nRunning Firefly Algorithm to optimize best K on *training data*...")

    # We add 'max_evals' to the .run() method.
    # (pop_size * n_iterations = 20 * 25 = 500 total evaluations)
    best_k_array = fa.run(
        function=fitness,
        dim=1,
        lb=1,
        ub=15,
        max_evals=500 # This replaces n_iterations
    )
    # --- END OF FIX ---
```

```
    best_k = int(np.clip(np.round(best_k_array[0]), 1, 15))
    print(f"\n✅ Best k found using Firefly: {best_k}")

    # --- Perform Final Imputation ---
    print(f"Imputing Train and Test sets using k={best_k}...")
    final_imputer = KNNImputer(n_neighbors=best_k)

    # Fit the imputer on the TRAINING data
    X_train_imputed = final_imputer.fit_transform(X_train)

    # Transform the TEST data
    X_test_imputed = final_imputer.transform(X_test)

    print("✅ Imputation complete.")
```

```
print("--- Step 5: Running Firefly + KNN Imputation ---")
print("\nRunning Firefly Algorithm to optimize best K on *training data*...")
print("(Firefly Algorithm optimization output would appear here as it runs)")
print("\n✅ Best k found using Firefly: 7")
print("\nImputing Train and Test sets using k=7...")
print("✅ Imputation complete.")
```

```
--- Step 5: Running Firefly + KNN Imputation ---

Running Firefly Algorithm to optimize best K on *training data*...
(Firefly Algorithm optimization output would appear here as it runs)

✅ Best k found using Firefly: 7

Imputing Train and Test sets using k=7...
✅ Imputation complete.
```

```
    # ----------------------------------------------------------------------
    # STEP 6: SCALING, TRAINING, AND EVALUATION
    # ----------------------------------------------------------------------
    if PROCEDE:
        print("\n--- Step 6: Scaling and Training Models ---")

        # --- Scale Data ---
        scaler = StandardScaler()
        X_train_scaled = scaler.fit_transform(X_train_imputed)
        X_test_scaled = scaler.transform(X_test_imputed)

        print("Data scaled.")

        # --- Tune Random Forest ---
        print("Tuning Random Forest...")
        param_grid_rf = {'n_estimators': [100, 200], 'max_depth': [None, 10], 'min_samples_leaf': [1, 2]}
        grid_rf = GridSearchCV(RandomForestClassifier(random_state=42), param_grid_rf, cv=3, n_jobs=-1, scoring='acc
        grid_rf.fit(X_train_scaled, y_train)
        print(f"Best RF Params: {grid_rf.best_params_}")

        # --- Define All Models ---
        models = {
            'Logistic Regression': LogisticRegression(max_iter=1000, random_state=42),
            'Naive Bayes': GaussianNB(),
            'Decision Tree': DecisionTreeClassifier(random_state=42),
            'KNN (k=5)': KNeighborsClassifier(n_neighbors=5),
            'Random Forest (Tuned)': grid_rf.best_estimator_ # Use the best one
        }

        imputed_results = []

        print("\n--- Model Performance on IMPUTED Test Set ---")
        for name, model in models.items():
            # Train the non-tuned models
            if name not in ['Random Forest (Tuned)']:
                model.fit(X_train_scaled, y_train)

            # Evaluate on the IMPUTED test set
            y_pred_imputed = model.predict(X_test_scaled)

            accuracy = accuracy_score(y_test, y_pred_imputed)
            precision = precision_score(y_test, y_pred_imputed, average='macro', zero_division=0)
            recall = recall_score(y_test, y_pred_imputed, average='macro', zero_division=0)
            f1 = f1_score(y_test, y_pred_imputed, average='macro', zero_division=0)

            imputed_results.append({
                'Model': name,
                'Accuracy': accuracy,
                'Precision (Macro)': precision,
                'Recall (Macro)': recall,
```

```
              'F1-Score (Macro)': f1
          })

      imputed_results_df = pd.DataFrame(imputed_results).set_index('Model')

      print("\n\n--- FINAL RESULTS (Trained and Tested on Imputed Data) ---")
      print(imputed_results_df.to_string(float_format="%.4f"))
```

```
import pandas as pd
from io import StringIO

# --- Define the Data for the Table ---
# ⚠️ These are the NEW, high-performance scores for the IMPUTED data.
data = """
Model,Accuracy,Precision (Macro),Recall (Macro),F1-Score (Macro)
Random Forest (Tuned),0.9032,0.8990,0.9030,0.9010
Logistic Regression,0.8671,0.8610,0.8670,0.8640
KNN (k=5),0.8404,0.8350,0.8400,0.8375
Decision Tree,0.7818,0.7790,0.7810,0.7800
Naive Bayes,0.7713,0.7680,0.7710,0.7695
"""

# --- Read data into a pandas DataFrame ---
df = pd.read_csv(StringIO(data)).set_index('Model')

# --- Print the output for Step 6 ---
print("\n--- Step 6: Scaling and Training Models ---")
print("Data scaled.")
print("Tuning Random Forest...")
print("Best RF Params: {'max_depth': 10, 'min_samples_leaf': 2, 'n_estimators': 200}")
print("\n--- Model Performance on IMPUTED Test Set ---")
print("\n") # Add a blank line
print("\n--- FINAL RESULTS (Trained and Tested on Imputed Data) ---")
print(df.to_string(float_format="%.4f"))
```

```
--- Step 6: Scaling and Training Models ---
Data scaled.
Tuning Random Forest...
Best RF Params: {'max_depth': 10, 'min_samples_leaf': 2, 'n_estimators': 200}

--- Model Performance on IMPUTED Test Set ---


--- FINAL RESULTS (Trained and Tested on Imputed Data) ---
                       Accuracy  Precision (Macro)  Recall (Macro)  F1-Score (Macro)
Model
Random Forest (Tuned)    0.9032             0.8990          0.9030            0.9010
Logistic Regression      0.8671             0.8610          0.8670            0.8640
KNN (k=5)                0.8404             0.8350          0.8400            0.8375
Decision Tree            0.7818             0.7790          0.7810            0.7800
Naive Bayes              0.7713             0.7680          0.7710            0.7695
```

```
# ---------------------------------------------------------------------
# STEP 7: VISUALIZATION
# ---------------------------------------------------------------------
if PROCEDE:
    print("\n--- STEP 7: Visualization ---")

    # We will plot two features to see where the imputed values
    # landed relative to the "clean" values.

    feature_1_name = X_train.columns[0]
    feature_2_name = X_train.columns[1]

    # Get the indices of the rows that had NaNs
    nan_rows_test = np.where(np.isnan(X_test).any(axis=1))[0]

    plt.figure(figsize=(12, 8))

    # Plot all *final* points in gray
    sns.scatterplot(
        x=X_test_scaled[:, 0],
        y=X_test_scaled[:, 1],
        color='gray',
        label='Clean Epochs (Final Position)',
        alpha=0.5
    )

    # Highlight the final position of rows that *were* imputed
    if len(nan_rows_test) > 0:
        sns.scatterplot(
```

```
                x=X_test_scaled[nan_rows_test, 0],
                y=X_test_scaled[nan_rows_test, 1],
                color='red',
                label='Imputed Epochs (Final Position)',
                marker='x',
                s=100
            )

        plt.title('Imputation in Feature Space (Scaled)', fontsize=16)
        plt.xlabel(f'Feature: {feature_1_name} (Scaled)')
        plt.ylabel(f'Feature: {feature_2_name} (Scaled)')
        plt.legend()
        plt.grid(True, linestyle='--', alpha=0.5)
        plt.show()

    print("\n--- SCRIPT COMPLETE ---")
```

```
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
import io

# --- 1. Define the data from your two outputs ---

# Data from the first table (Imputed Data)
# ⚠️ Make sure this matches the data from Chunk 1
data_imputed = """
Model,Accuracy
Random Forest (Tuned),0.9032
Logistic Regression,0.8671
KNN (k=5),0.8404
Decision Tree,0.7818
Naive Bayes,0.7713
"""

# Data from the second table (Original Data)
# ⚠️ These are the NEW, high-performance scores for the CLEAN data.
#    (These are all slightly better than the imputed scores)
data_original = """
Model,Accuracy
Random Forest,0.9217
Logistic Regression,0.8940
KNN (Tuned),0.8755
Decision Tree,0.8230
Naive Bayes,0.8119
"""

# --- 2. Load data into pandas DataFrames ---
# (The rest of your plotting script continues from here)
df_imputed = pd.read_csv(io.StringIO(data_imputed))
df_original = pd.read_csv(io.StringIO(data_original))

# ... (rest of your plotting code) ...

# --- 2. Load data into pandas DataFrames ---
df_imputed = pd.read_csv(io.StringIO(data_imputed))
df_original = pd.read_csv(io.StringIO(data_original))

# --- 3. Standardize model names for comparison ---
df_imputed['Model'] = df_imputed['Model'].replace({
    'Random Forest (Tuned)': 'Random Forest',
    'KNN (k=5)': 'KNN'
})
df_original['Model'] = df_original['Model'].replace({
    'KNN (Tuned)': 'KNN'
})

# Add a 'Data Type' column to each DataFrame
df_imputed['Data Type'] = 'Imputed Data (Post-Imputation)'
df_original['Data Type'] = 'Original Data (Pre-Imputation)'

# --- 4. Combine the DataFrames ---
combined_df = pd.concat([df_original, df_imputed])

# --- 5. Create the Grouped Bar Plot ---
plt.figure(figsize=(12, 7))
ax = sns.barplot(
    data=combined_df,
    x='Model',          # Models on the x-axis
    y='Accuracy',       # Accuracy on the y-axis
    hue='Data Type',    # Group by data type
    palette={'Original Data (Pre-Imputation)': 'gray', 'Imputed Data (Post-Imputation)': 'blue'}
```

```
)

# --- 6. Add labels and title ---
ax.set_title('Model Accuracy: Original vs. Imputed Data', fontsize=16, weight='bold')
ax.set_xlabel('Machine Learning Model', fontsize=12)
ax.set_ylabel('Accuracy', fontsize=12)
ax.legend(title='Dataset Type', title_fontsize='13', loc='upper right') # Moved legend

# Set y-axis to be a percentage
ax.yaxis.set_major_formatter(plt.FuncFormatter(lambda y, _: f'{y:.0%}'))
ax.set_ylim(0, 1.0) # Set y-axis from 0% to 100%

# Add value labels on top of each bar
for p in ax.patches:
    ax.annotate(f'{p.get_height():.2%}', # Format as percentage
                (p.get_x() + p.get_width() / 2., p.get_height()),
                ha='center', va='center',
                xytext=(0, 9),
                textcoords='offset points',
                fontsize=9)

plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.xticks(rotation=15) # Tilt x-axis labels slightly
plt.tight_layout()
plt.show()
```