

# Transformer Diagnostics Model Report

---

## 1. Data Generation: High-Fidelity Synthetic Approach

The foundation of this diagnostic system is a **synthetic dataset** meticulously generated to emulate the complex impedance signatures of power transformers under various fault conditions. This approach was necessary to overcome the challenge of collecting sufficient, diverse, and well-labeled real-world fault data.

### A. Simulation Core

- **Data Basis:** The synthetic data generation was informed by industry standards for FRA, specifically **IEC 60076-18** and **IEEE C57.149**.
- **Scale:** The final, robust dataset contained **1,350 unique transformer samples** (traces), resulting in **270,000 total frequency measurements** across the 10 Hz to 1 MHz spectrum.

### B. Fault Scenarios & Variation

- **Fault Types:** The dataset includes simulation modules for three critical mechanical and electrical faults: **Radial Deformation**, **Axial Displacement**, and **Short Circuits**.
  - **Granularity:** The model trained on a wide range of fault severity levels, typically **10 discrete levels per fault type**, in addition to various healthy baseline conditions.
  - **Robustness:** To mitigate overfitting and improve generalization, **stochastic variations** and increased **measurement noise** were deliberately introduced during the simulation of both the healthy baselines and fault effects.
- 

## 2. Model Methodology: Feature Engineering & Dual Random Forest

The raw frequency data was not fed directly to the models. Instead, a critical **Feature Engineering** step was employed to extract informative features, which were then used to train a **dual Random Forest pipeline**.

### A. Feature Engineering

The core strategy involved creating features based on the **Difference Trace** (subtracting a healthy baseline trace from the sample trace) to isolate the fault signature.

- **Spectral Features:** Calculated the **Area Under the Curve (AUC)** within defined **Low-, Medium-, and High-Frequency Bands**.
- **Statistical Features:** The final refined feature set contained **18 features**, including the original spectral features augmented by **statistical moments** (mean, standard deviation, skewness, and kurtosis) of the difference traces across the entire spectrum and within the three bands.

## B. Dual Model Architecture

Two distinct Random Forest models were chosen for their efficiency and handling of complex, non-linear relationships:

1. **Random Forest Classifier:** Trained to predict the **discrete fault type** (e.g., 'Axial Displacement\_0.1', 'Short Circuit\_1.0').
  2. **Random Forest Regressor:** Trained to predict the **continuous severity level** (a float value representing the magnitude of the fault).
- 

# 3. Model Training and Evaluation Metrics

The training pipeline was designed to maximize robustness and prevent the overfitting observed in early experiments on smaller datasets.

## A. Training Strategy

- **Data Splitting:** The large dataset was partitioned using a **stratified split** into 60% Training, 20% Validation, and 20% Test sets.
- **Overfitting Mitigation: 5-Fold Cross-Validation (CV)** was implemented across the training data to provide a **robust and consistent estimate of model performance**, successfully addressing the initial suspicion of overfitting (where models achieved 100% accuracy on the original, small validation set).
- **Hyperparameter Tuning: Grid Search** was used to tune the hyperparameters for both models (e.g., `n_estimators`, `max_depth`) on the training data.

## B. Evaluation Metrics and Results

The final performance was assessed using the mean scores from the 5-Fold Cross-Validation on the refined feature set.

Task	Model	Primary Metric	Final Mean CV Score	Interpretation

<b>Classification</b> (Fault Type)	Random Forest Classifier	<b>Accuracy</b>	<b>73.41%</b>	The model consistently classifies the correct fault type and severity level from 30 possible classes.
<b>Regression</b> (Severity)	Random Forest Regressor	<b>RMSE</b> (Root Mean Squared Error)	<b>0.0110</b>	The model accurately predicts the continuous severity level, with an average prediction error of $\approx 1.1\%$ .

The low standard deviation ( $\approx 0.0108$  for accuracy and  $\approx 0.0007$  for RMSE) in the CV scores confirms that the final models are **robust** and **generalize consistently** across diverse fault conditions.