

## Exploratory Data Analysis For CAMPSS Dataset

### Importing important Libraries.

```
# 3/85 -> Block number to track process during execution

# All Libraries are retrospectively pasted here for code clearance
# Standard Library
import warnings
# Data Manipulation & Visualization
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px

# Machine Learning & Preprocessing
import sklearn
from sklearn.ensemble import RandomForestRegressor # selected as an alternative to XGBoost
from sklearn.model_selection import train_test_split, TimeSeriesSplit
from sklearn.preprocessing import StandardScaler, RobustScaler
from sklearn.metrics import mean_squared_error as MSE,
    mean_squared_error as MSE,
    mean_absolute_percentage_error as MAPE

# Statistical Tools
from scipy.stats import zscore, gaussian_kde

# Imbalanced Data Handling
from imblearn.under_sampling import RandomUnderSampler
from sklearn.model_selection import GroupShuffleSplit

# Gradient Boosting
import xgboost as xgb

# Jupyter-specific
%matplotlib inline

# Suppress Warnings
warnings.filterwarnings('ignore')

print("Library Load Successful")
```

### Reading all the "train\_FDOO[i].txt" files as Dataframes.

```
# 5/85
train1=pd.read_csv('train_FDO01.txt',sep='*',header=None)
train2=pd.read_csv('train_FDO02.txt',sep='*',header=None)
train3=pd.read_csv('train_FDO03.txt',sep='*',header=None)
train4=pd.read_csv('train_FDO04.txt',sep='*',header=None)

test1 = pd.read_csv('test_FDO01.txt', sep='*', header=None)
test2 = pd.read_csv('test_FDO02.txt', sep='*', header=None)
test3 = pd.read_csv('test_FDO03.txt', sep='*', header=None)
test4 = pd.read_csv('test_FDO04.txt', sep='*', header=None)

rul1 = pd.read_csv('RUL_FDO01.txt', header=None, names=['RUL'])
rul2 = pd.read_csv('RUL_FDO02.txt', header=None, names=['RUL'])
rul3 = pd.read_csv('RUL_FDO03.txt', header=None, names=['RUL'])
rul4 = pd.read_csv('RUL_FDO04.txt', header=None, names=['RUL'])

print("Reading .txt to dataframes Successful")

Reading .txt to dataframes Successful
```

### Getting the shapes of all the "train\_FDOO[i].txt" files.

```
# 7/85
print("Shape of 'train_FDO01.txt' is: ",train1.shape)
print("Shape of 'train_FDO02.txt' is: ",train2.shape)
print("Shape of 'train_FDO03.txt' is: ",train3.shape)
print("Shape of 'train_FDO04.txt' is: ",train4.shape)
print("-----")
print("Shape of 'test_FDO01.txt' is: ",test1.shape)
print("Shape of 'test_FDO02.txt' is: ",test2.shape)
print("Shape of 'test_FDO03.txt' is: ",test3.shape)
print("Shape of 'test_FDO04.txt' is: ",test4.shape)
print("-----")
print("Shape of 'RUL_FDO01.txt' is: ",rul1.shape)
print("Shape of 'RUL_FDO02.txt' is: ",rul2.shape)
print("Shape of 'RUL_FDO03.txt' is: ",rul3.shape)
print("Shape of 'RUL_FDO04.txt' is: ",rul4.shape)
print("-----")
print("Shape printing Successful")
```

Shape of 'train\_FDO01.txt' is: (20631, 26)

Shape of 'train\_FDO02.txt' is: (53759, 26)

Shape of 'train\_FDO03.txt' is: (13996, 26)

Shape of 'train\_FDO04.txt' is: (61249, 26)

-----

Shape of 'test\_FDO01.txt' is: (13996, 26)

Shape of 'test\_FDO02.txt' is: (33991, 26)

Shape of 'test\_FDO03.txt' is: (16596, 26)

Shape of 'test\_FDO04.txt' is: (31121, 26)

-----

Shape of 'RUL\_FDO01.txt' is: (13996, 1)

Shape of 'RUL\_FDO02.txt' is: (2529, 1)

Shape of 'RUL\_FDO03.txt' is: (108, 1)

Shape of 'RUL\_FDO04.txt' is: (248, 1)

-----

Shape printing Successful

### Displaying first 5 rows of "train\_FDO01.txt"

```
# 9/85
#train1
print("First 5 rows of 'train_FDO01.txt' are: ")
train1.head()

First 5 rows of 'train_FDO01.txt' are:
   0   1   2   3   4   5   6   7   8   9   ...   16   17   18   19   20   21   22   23   24   25
0  1  100.0  516.67  641.82  1589.70  1400.60  14.62 ... 521.66  2388.02  8138.62  8.4195  0.03  392  2388  100.0  39.06  23.4190
1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
2  1  3  0.0043  0.0003  100.0  516.67  642.35  1587.99  1404.20  14.62 ... 522.42  2388.03  8133.23  8.4178  0.03  390  2388  100.0  38.95  23.3442
3  1  4  0.0007  0.0000  100.0  516.67  642.35  1582.79  1401.87  14.62 ... 522.88  2388.08  8133.83  8.3692  0.03  392  2388  100.0  38.88  23.3739
4  1  5  0.0019  -0.0002  100.0  516.67  642.37  1582.85  1406.22  14.62 ... 522.19  2388.04  8133.80  8.4294  0.03  393  2388  100.0  38.90  23.4044
5 rows < 26 columns
```

### Providing the Feature names.

```
# 11/85
sensor_names = ['unit',
                 'cycle',
                 'opset1',
                 'opset2',
                 'opset3',
                 'opset4',
                 'Fan_Inlet_Temp',
                 'LPC_Outlet_Temp',
                 'HPC_Outlet_Temp',
                 'LPT_Outlet_Temp',
                 'Fan_Inlet_Press',
                 'Bypass_Duct',
                 'HPC_Outlet_Press',
                 'HPC_Outlet_Speed',
                 'Physical_Core_Speed',
                 'Engine_Press_Ratio',
                 'Fuel_Flow_Rate',
                 'HPC_Outlet_Static_Press',
                 'Fuel_Air_Ratio',
                 'Bleed_Enthalpy',
                 'Req_Fan_Speed',
                 'Req_Conv_Speed',
                 'HPT_Cool_Air_F',
                 'LPT_Cool_Air_F'
                ]
# Assigning the feature names to the corresponding Dataframe.
datasets = [train1, train2, train3, train4, test1, test2, test3, test4]
for df in datasets:
    df.columns = sensor_names
print("Columns renamed. Sample train1 columns:", train1.columns.tolist()[1:5])
print("train1 head:\n", train1.head(2))

train1.head()
unit  cycle  opset1  opset2  opset3  Fan_Inlet_Temp  LPC_Outlet_Temp  HPC_Outlet_Temp  LPT_Outlet_Temp  Fan_Inlet_Press  Bypass_Duct  HPC_Outlet_Press  Fuel_Flow_Rate  HPC_Outlet_Static_Press  Fuel_Air_Ratio  Bleed_Enthalpy  Req_Fan_Speed  Req_Conv_Speed  HPT_Cool_Air_F  LPT_Cool_Air_F
0  1  100.0  516.67  641.82  1589.70  1400.60  14.62 ... 521.66  2388.02  8138.62  8.4195  0.03  392  2388  100.0  39.06  23.4190
1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
2  1  3  0.0043  0.0003  100.0  516.67  642.35  1587.99  1404.20  14.62 ... 522.42  2388.03  8133.23  8.4178  0.03  390  2388  100.0  38.95  23.3442
3  1  4  0.0007  0.0000  100.0  516.67  642.35  1582.79  1401.87  14.62 ... 522.88  2388.08  8133.83  8.3692  0.03  392  2388  100.0  38.88  23.3739
4  1  5  0.0019  -0.0002  100.0  516.67  642.37  1582.85  1406.22  14.62 ... 522.19  2388.04  8133.80  8.4294  0.03  393  2388  100.0  38.90  23.4044
5  2  3  0.0007  0.0000  100.0  516.67  642.37  1582.85  1406.22  14.62 ... 522.42  2388.03  8133.80  8.4294  0.03  393  2388  100.0  38.90  23.4044
6  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
7  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
8  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
9  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
10  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
11  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
12  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
13  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
14  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
15  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
16  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
17  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
18  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
19  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
20  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
21  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
22  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
23  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
24  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
25  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
26  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
27  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100.0  39.00  23.4236
28  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  2  0.0019  -0.0003  100.0  516.67  642.19  1591.82  1403.14  14.62 ... 522.28  2388.07  8131.49  8.4318  0.03  392  2388  100
```

```

28 Fuel_Air_Ratio 20831 non-null float64
21 Bleed_Enthalpy 20831 non-null int64
22 Req_Fan_Speed 20831 non-null int64
23 Req_Fan_Speed 20831 non-null int64
24 HPT_Cool_Air_F 20831 non-null float64
25 LPT_Cool_Air_F 20831 non-null float64
dtype: object
memory usage: 4.1 MB
=====
```

Null Values of all the train Dataframes.

```

# 17/85
print("Null Values of 'train_F0001.txt' file are: ")
print("Null Values of 'train_F0002.txt' file are: ")
print("Null Values of 'train_F0003.txt' file are: ")
print("Null Values of 'train_F0004.txt' file are: ")
print("Null Values of 'train_F0005.txt' file are: ")
```

Null Values of 'train\_F0001.txt' file are: 0

```

unit 0
cycle 0
opset1 0
opset2 0
opset3 0
Fan_Inlet_Temp 0
LPC_Outlet_Temp 0
HPC_Outlet_Temp 0
LPT_Outlet_Temp 0
Fan_Inlet_Press 0
Bypass_Duct 0
HPC_Outlet_Press 0
Physical_Fan_Speed 0
Physical_Core_Speed 0
Engine_Press_Ratio 0
HPC_Outlet_Static_Press 0
Fuel_Flow_Ratio 0
Corrected_Fan_Speed 0
Corrected_Core_Speed 0
Bypass_Ratio 0
Fuel_Air_Ratio 0
Bleed_Enthalpy 0
Req_Fan_Speed 0
Req_Core_Speed 0
HPT_Cool_Air_F 0
LPT_Cool_Air_F 0
dtype: int64
```

Null Values of 'train\_F0002.txt' file are: 0

```

unit 0
cycle 0
opset1 0
opset2 0
opset3 0
Fan_Inlet_Temp 0
LPC_Outlet_Temp 0
HPC_Outlet_Temp 0
LPT_Outlet_Temp 0
Fan_Inlet_Press 0
Bypass_Duct 0
HPC_Outlet_Press 0
Physical_Fan_Speed 0
Physical_Core_Speed 0
Engine_Press_Ratio 0
HPC_Outlet_Static_Press 0
Fuel_Flow_Ratio 0
Corrected_Fan_Speed 0
Corrected_Core_Speed 0
Bypass_Ratio 0
Fuel_Air_Ratio 0
Bleed_Enthalpy 0
Req_Fan_Speed 0
Req_Core_Speed 0
HPT_Cool_Air_F 0
LPT_Cool_Air_F 0
dtype: int64
```

Checking for the duplicated values in all the train Dataframes.

```

# 19/85
print("Duplicated values in 'train1' Dataframe are: {train1.duplicated().sum()}")
print("Duplicated values in 'train2' Dataframe are: {train2.duplicated().sum()}")
print("Duplicated values in 'train3' Dataframe are: {train3.duplicated().sum()}")
print("Duplicated values in 'train4' Dataframe are: {train4.duplicated().sum()}")
print("Duplicated values in 'train5' Dataframe are: {train5.duplicated().sum()}")
```

Duplicated values in 'train1' Dataframe are: 0

Duplicated values in 'train2' Dataframe are: 0

Duplicated values in 'train3' Dataframe are: 0

Duplicated values in 'train4' Dataframe are: 0

Describing the Dataframes.

# 21/85

```

#train1
print("Describing the 'train1' Dataframe:")
train1.describe(),T
```

Describing the 'train1' Dataframe:

	count	mean	std	min	25%	50%	75%	max
unit	20631.0	51.506598	2.022763e+01	1.0000	26.0000	62.0000	77.0000	100.0000
cycle	20631.0	108.807062	6.88009e+01	1.0000	52.0000	104.0000	156.0000	362.0000
opset1	20631.0	-0.00009	2.18713e-03	-0.0087	-0.0015	0.0000	0.0015	0.0087
opset2	20631.0	0.00002	2.90021e-04	-0.0006	-0.0002	0.0000	0.0003	0.0006
opset3	20631.0	100.000000	2.90002e+00	100.0000	100.0000	100.0000	100.0000	100.0000
Fan_Inlet_Temp	20631.0	518.670000	0.00000e+00	518.6700	518.6700	518.6700	518.6700	518.6700
LPC_Outlet_Temp	20631.0	642.689034	5.00053e+00	641.2100	642.3250	642.6400	643.0000	644.5300
HPC_Outlet_Temp	20631.0	1590.523119	6.13150e+00	1571.0400	1582.2600	1591.1000	1594.3800	1616.9100
LPT_Outlet_Temp	20631.0	1408.933782	9.00005e+00	1382.2500	1402.3600	1408.0400	1414.5500	1441.4900
Fan_Inlet_Press	20631.0	14.620000	5.329200e-15	14.6200	14.6200	14.6200	14.6200	14.6200
Bypass_Duct	20631.0	21.609803	1.388985e-03	21.6000	21.6100	21.6100	21.6100	21.6100
HPC_Outlet_Press	20631.0	553.367711	8.850926e-01	549.8600	552.8100	553.4400	554.0100	556.0600
Physical_Fan_Speed	20631.0	2388.096952	7.098548e-02	2387.9000	2386.0500	2386.0000	2388.1400	2388.5600
Physical_Core_Speed	20631.0	9065.242941	2.028288e+01	9021.7300	9053.1000	9060.6600	9069.4200	9244.5900
Engine_Press_Ratio	20631.0	1.300000	0.00000e+00	1.3000	1.3000	1.3000	1.3000	1.3000
HPC_Outlet_Static_Press	20631.0	47.541168	2.670747e-01	46.8500	47.3500	47.5100	47.7000	48.5300
Fuel_Flow_Ratio	20631.0	521.143470	3.735534e-01	518.6900	520.9600	521.4800	521.9500	523.5800
Corrected_Fan_Speed	20631.0	2388.096952	7.191952e-02	2387.8800	2388.0400	2386.0000	2388.1400	2388.5600
Corrected_Core_Speed	20631.0	8143.752722	1.90718e+01	8099.8400	8130.2450	8140.5400	8148.8100	8293.7200
Bypass_Ratio	20631.0	8.442146	3.750504e-02	8.3249	8.4149	8.4389	8.4656	8.5848
Fuel_Air_Rate	20631.0	0.030000	3.498531e-18	0.0300	0.0300	0.0300	0.0300	0.0300
Bleed_Enthalpy	20631.0	392.210654	1.548763e+00	388.0000	392.0000	393.0000	394.0000	400.0000
Req_Fan_Speed	20631.0	2388.000000	0.00000e+00	2388.0000	2388.0000	2388.0000	2388.0000	2388.0000
Req_Core_Speed	20631.0	100.000000	0.00000e+00	100.0000	100.0000	100.0000	100.0000	100.0000
HPT_Cool_Air_F	20631.0	38.816271	1.807464e-01	38.1400	38.7000	38.8300	38.9500	39.4300
LPT_Cool_Air_F	20631.0	23.289705	1.082509e-01	22.8942	23.2218	23.2979	23.3668	23.5184

# 22/85

```

#train2
print("Describing the 'train2' Dataframe:")
train2.describe(),T
```

Describing the 'train2' Dataframe:

	count	mean	std	min	25%	50%	75%	max
unit	53759.0	131.082981	74.463862	1.0000	68.0000	131.0000	195.0000	260.0000
cycle	53759.0	109.154746	69.180569	1.0000	52.0000	104.0000	157.0000	378.0000
opset1	53759.0	23.998407	14.747376	0.0000	10.0040	25.0013	41.9980	42.0080
opset2	53759.0	0.572056	0.310016	0.0000	0.2507	0.7000	0.8400	0.8420
opset3	53759.0	94.040020	14.237738	60.0000	100.0000	100.0000	100.0000	100.0000
Fan_Inlet_Temp	53759.0	472.910207	26.389707	445.0000	445.0000	462.5400	491.1900	518.6700
LPC_Outlet_Temp	53759.0	579.672999	37.280399	535.5300	549.5700	556.9800	607.3400	644.5200
HPC_Outlet_Temp	53759.0	1419.971013	105.946341	124.7300	352.7600	339.1800	449.3700	1612.8800
LPT_Outlet_Temp	53759.0	1205.442024	119.123428	102.7700	1123.6500	1138.8900	1306.8500	1439.2300
Fan_Inlet_Press	53759.0	8.031686	3.613830	3.9100	3.9100	7.0500	10.5200	14.6200
Bypass_Duct	53759.0	11.600746	5.431802	5.7100	5.7200	9.0300	15.4900	21.6100
HPC_Outlet_Press	53759.0	282.606787	146.005300	136.8000	139.9300	194.6600	394.0800	555.8200
Physical_Fan_Speed	53759.0	2228.078188	145.209816	191.7700	221.8800	223.0700	2323.9600	2388.3900
Physical_Core_Speed	53759.0	8252.209837	335.812019	7885.5600	8321.6600	8361.2000	8778.0300	9215.6600
Engine_Press_Ratio	53759.0	1.094962	0.127400	0.9300	1.0200	1.0200	1.2600	1.3000
HPC_Outlet_Static_Press	53759.0	42.985172	3.232372	36.2300	41.9100	42.3900	45.3500	48.5100
Fuel_Flow_Ratio	53759.0	266.069034	137.659507	128.1200	131.5200	183.2000	371.2600	523.3700
Corrected_Fan_Speed	53759.0	2334.57253	128.068271	2027.6100	2387.9000	2388.0800	2388.1700	2390.4800
Corrected_Core_Speed	53759.0	6006.597682	84.837950	7848.3600	8062.1400	8082.5400	8172.1950	8268.5000
Bypass_Ratio	53759.0	9.326954	0.749330	8.3357	8.6774	9.3100	9.3869	11.0669
Fuel_Air_Rate	53759.0	0.023326	0.004711	0.0200	0.0200	0.0200	0.0300	0.0300
Bleed_Enthalpy	53759.0	348.309511	27.754516	303.0000	331.0000	335.0000	369.0000	399.0000
Req_Fan_Speed	53759.0	2228.060358	145.327980	1915.0000	2212.0000	2223.0000	2324.0000	2388.0000
Req_Core_Speed	53759.0	97.750638	5.364067	84.9300	100.0000	100.0000	100.0000	100.0000
HPT_Cool_Air_F	53759.0	20.789296	9.869333	10.1900	10.9100	14.8800	28.4700	39.3400
LPT_Cool_Air_F	53759.0	12.						

Describing the 'train1' Dataframe:								
	count	mean	std	min	25%	50%	75%	max
unit	61249.0	124.325181	71.995350	1.0000	60.0000	126.0000	185.0000	249.0000
cycle	61249.0	134.311417	88.783389	1.0000	62.0000	123.0000	191.0000	543.0000
opset1	61249.0	23.999823	14.780722	0.0000	10.0046	25.0014	41.9981	42.0000
opset2	61249.0	0.571347	0.310703	0.0000	0.2507	0.7000	0.8400	0.8420
opset3	61249.0	94.031676	14.251954	60.0000	100.0000	100.0000	100.0000	100.0000
Fan_Inlet_Temp	61249.0	472.882435	26.436832	445.0000	445.0000	462.5400	491.1900	518.6700
LPC_Outlet_Temp	61249.0	579.420056	37.342647	535.4800	549.3300	555.7400	607.0700	644.4200
HPC_Outlet_Temp	61249.0	1417.896000	106.167598	1242.6700	1350.5500	1367.6800	1497.4200	1613.0000
LPT_Outlet_Temp	61249.0	1201.915359	119.327591	1024.4200	1119.4900	1136.9200	1302.6200	1440.7700
Fan_Inlet_Press	61249.0	8.031628	3.622872	3.9100	3.9100	7.0500	10.5200	14.6200
Bypass_Duct	61249.0	11.589457	5.444017	5.6700	5.7200	9.0300	15.4800	21.6100
HPC_Outlet_Press	61249.0	283.328633	146.880210	136.1700	142.9200	194.9600	394.2800	570.8100
Physical_Fan_Speed	61249.0	2228.680034	145.348243	1914.7200	2211.9500	2223.0700	2233.9300	2388.6400
Physical_Core_Speed	61249.0	8524.673301	336.927547	7984.5100	8320.5900	8362.7600	8777.2500	9196.8100
Engine_Press_Ratio	61249.0	1.064445	0.127681	0.9500	1.0200	1.0300	1.2600	1.3200
HPC_Outlet_Static_Press	61249.0	42.874529	3.243482	36.0400	41.7600	45.2200	48.3600	
Fuel_Flow_Ratio	61249.0	266.735665	138.473490	128.3100	134.5200	183.4500	371.4000	537.4900
Corrected_Core_Speed	61249.0	2334.427590	128.197859	2027.5700	2387.9100	2388.0600	2388.1700	2390.4900
Corrected_Fan_Speed	61249.0	8067.818182	85.670543	7845.7800	8062.6300	8083.8100	8128.3500	8261.6500
Bypass_Ratio	61249.0	9.285604	0.750374	8.1757	8.6480	9.2556	9.3658	11.0663
Fuel_Air_Ratio	61249.0	0.023252	0.004685	0.0200	0.0200	0.0200	0.0300	0.0300
Bleed_Enthalpy	61249.0	347.760029	27.808283	302.0000	330.0000	334.0000	368.0000	399.0000
Req_Fan_Speed	61249.0	2228.613283	145.474241	1915.0000	2212.0000	2223.0000	2324.0000	2388.0000
Req_Core_Speed	61249.0	97.751396	5.369424	84.9300	100.0000	100.0000	100.0000	100.0000
HPT_Cool_Air_F	61249.0	20.884333	0.936306	10.1600	10.9400	14.9300	28.5600	39.8900
LPT_Cool_Air_F	61249.0	12.518995	5.962697	6.0843	6.5561	8.9601	17.1355	23.8882

↳ Unique categories of each feature in train Dataframes.

```
# 26/85
print("Unique categories of each feature in train1 Dataframe: \n{len(train1.columns)}\n\n")
print("-----")
print("Unique categories of each feature in train2 Dataframe: \n{len(train2.columns)}\n\n")
print("-----")
print("Unique categories of each feature in train3 Dataframe: \n{len(train3.columns)}\n\n")
print("-----")
print("Unique categories of each feature in train4 Dataframe: \n{len(train4.columns)}\n\n")
```

Unique categories of each feature in train1 Dataframe: 26

unit	108
cycle	362
opset1	158
opset2	13
opset3	1
Fan_Inlet_Temp	1
LPC_Outlet_Temp	310
HPC_Outlet_Temp	3812
LPT_Outlet_Temp	4051
Fan_Inlet_Press	1
Bypass_Duct	2
HPC_Outlet_Press	513
Physical_Fan_Speed	53
Physical_Core_Speed	6463
Engine_Press_Ratio	1
HPC_Outlet_Static_Press	159
Fuel_Flow_Ratio	427
Corrected_Fan_Speed	56
Corrected_Core_Speed	609
Bypass_Ratio	1918
Fuel_Air_Ratio	1
Bleed_Enthalpy	11
Req_Fan_Speed	1
HPT_Cool_Air_F	138
LPT_Cool_Air_F	4745
	dtype: int64

Unique categories of each feature in train2 Dataframe: 26

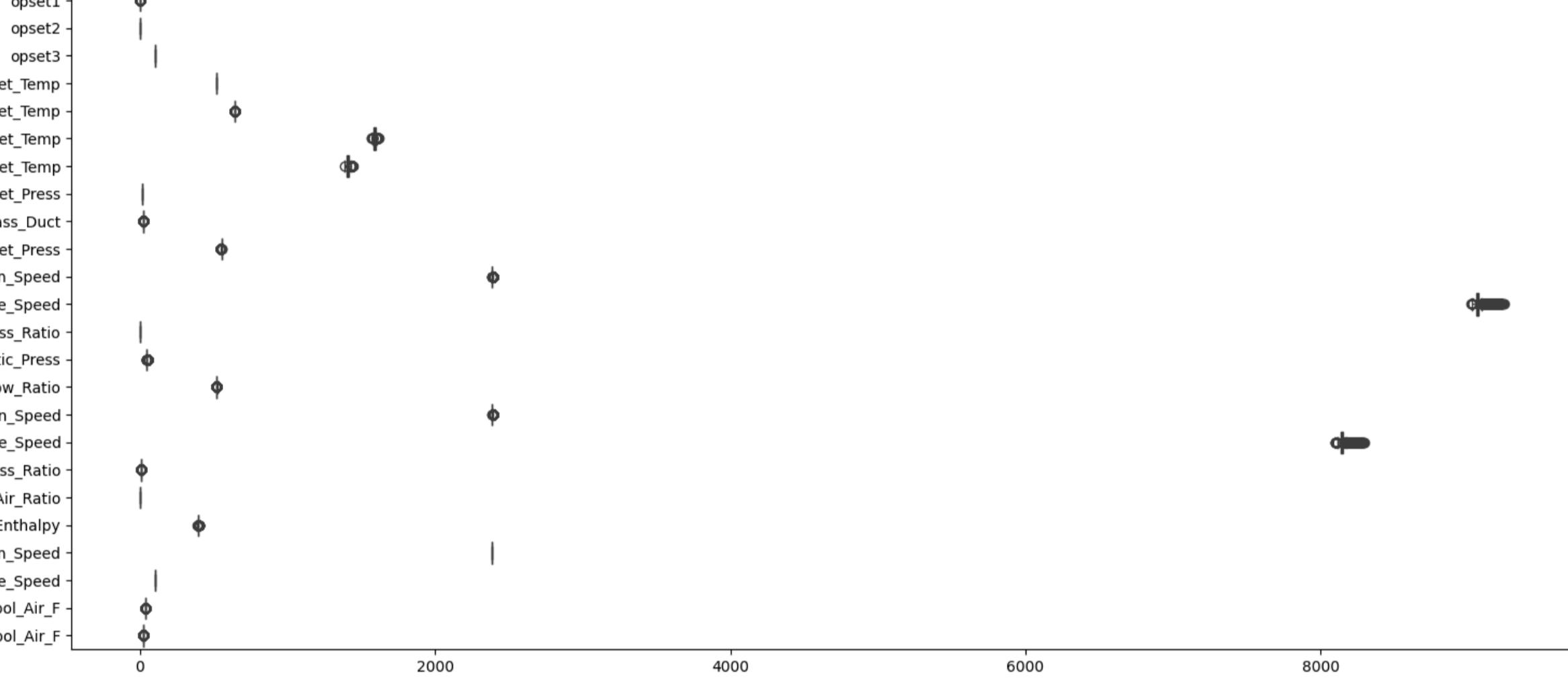
unit	268
cycle	378
opset1	536
opset2	105
opset3	2
Fan_Inlet_Temp	6
LPC_Outlet_Temp	1598
HPC_Outlet_Temp	11095
LPT_Outlet_Temp	15411
Fan_Inlet_Press	6
Bypass_Duct	14
HPC_Outlet_Press	2067
Physical_Fan_Speed	897
Physical_Core_Speed	2244
Engine_Press_Ratio	9
HPC_Outlet_Static_Press	688
Fuel_Flow_Ratio	1672
Corrected_Fan_Speed	514
Corrected_Core_Speed	14985
Bypass_Ratio	8464
Fuel_Air_Ratio	2
Bleed_Enthalpy	53
Req_Fan_Speed	6
Req_Core_Speed	2
HPT_Cool_Air_F	510

Outlier Analysis

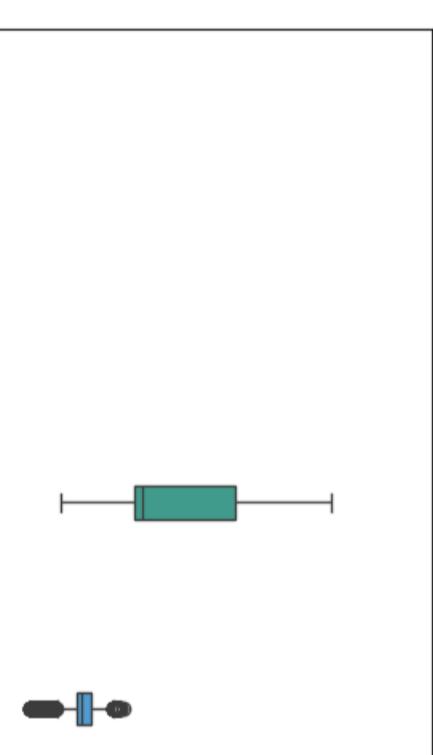
```
# 28/85
# Boxplot outlier analysis for all sensors
def plot_outliers(df, title):
    plt.figure(figsize=(18, 8))
    sns.boxplot(x=df, orient="h")
    plt.title(title)
    plt.show()

plot_outliers(train1.drop(["unit", "cycle"], axis=1), "train_FD001 Sensor Outlier Analysis")
plot_outliers(train2.drop(["unit", "cycle"], axis=1), "train_FD002 Sensor Outlier Analysis")
plot_outliers(train3.drop(["unit", "cycle"], axis=1), "train_FD003 Sensor Outlier Analysis")
plot_outliers(train4.drop(["unit", "cycle"], axis=1), "train_FD004 Sensor Outlier Analysis")
```

train\_FD001 Sensor Outlier Analysis



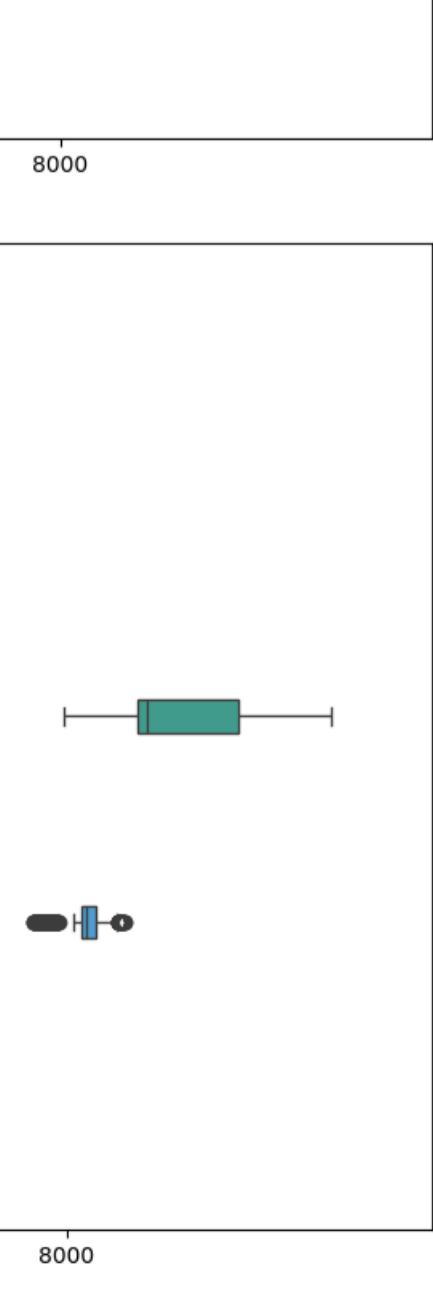
train\_FD002 Sensor Outlier Analysis



train\_FD003 Sensor Outlier Analysis



train\_FD004 Sensor Outlier Analysis



Z-Score

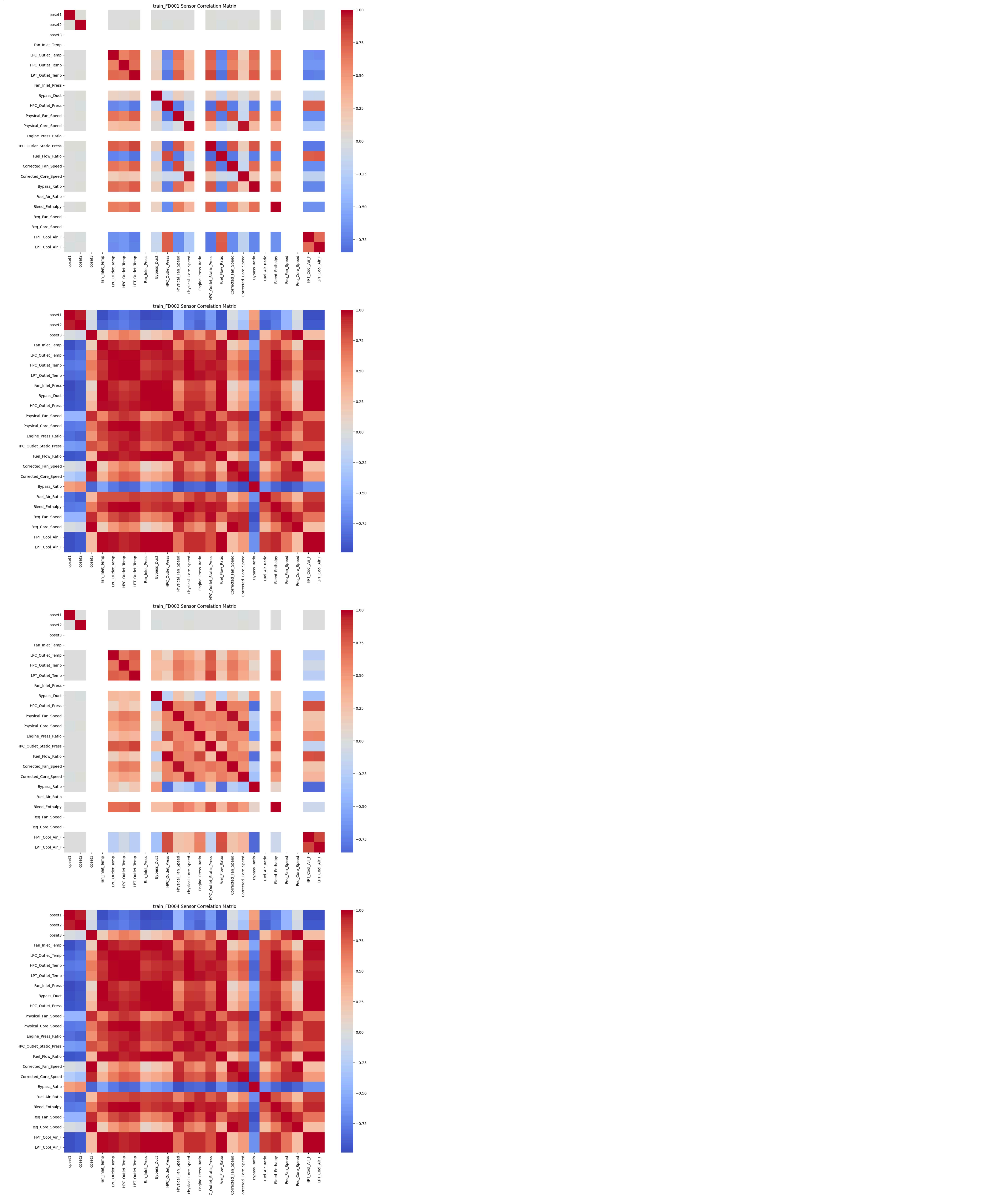
```
# 30/85
# Z-score > 3 or < -3 are considered outliers"
def zscore_outlier_summary(df, name):
    df_numeric = df.drop(["unit", "cycle"], axis=1)
    z_scores = np.abs(zscore(df_numeric))
```

## Correlation Matrix and Heat Map

```
# 32/85

def plot_corr_heatmap(df, title):
    df_numeric = df.drop(["unit", "cycle"], axis=1)
    corr = df_numeric.corr()
    plt.figure(figsize=(16, 12))
    sns.heatmap(corr, annot=False, cmap="coolwarm", center=0)
    plt.title(title)
    plt.show()

plot_corr_heatmap(train1, "train_FD001 Sensor Correlation Matrix")
plot_corr_heatmap(train2, "train_FD002 Sensor Correlation Matrix")
plot_corr_heatmap(train3, "train_FD003 Sensor Correlation Matrix")
plot_corr_heatmap(train4, "train_FD004 Sensor Correlation Matrix")
```

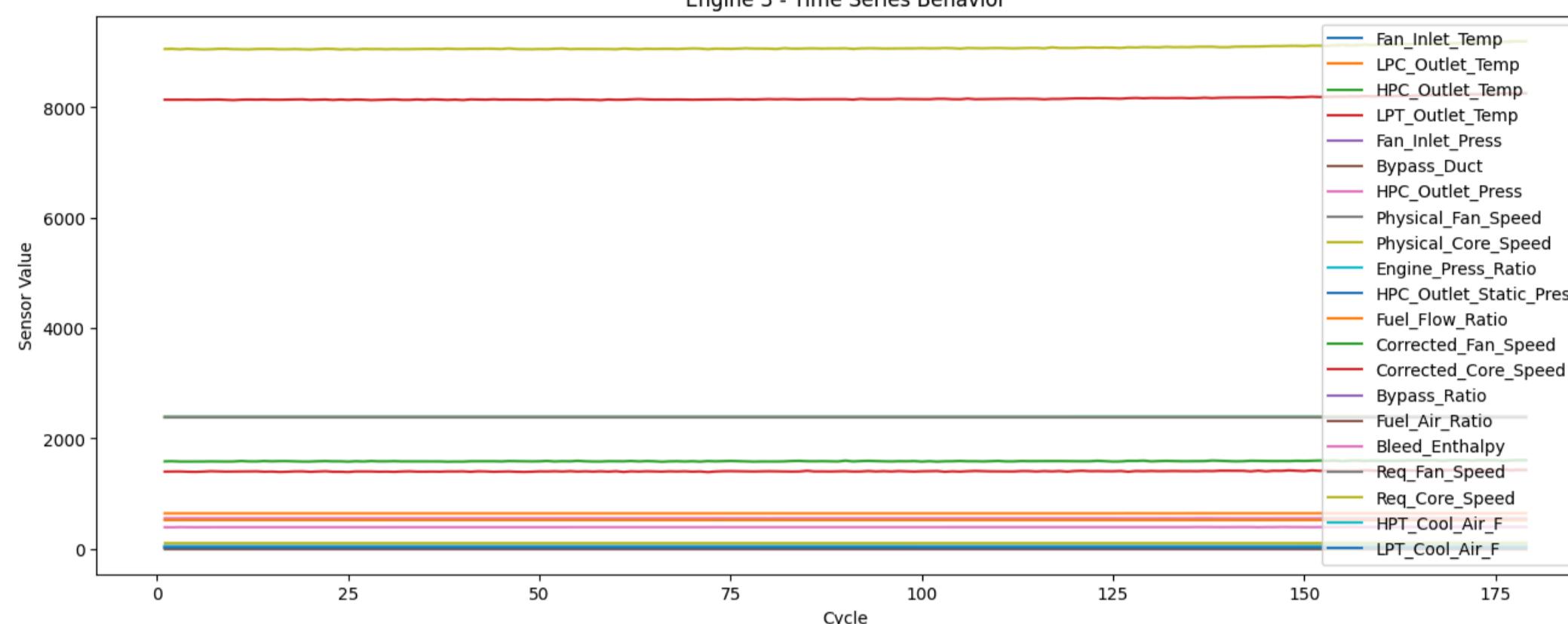
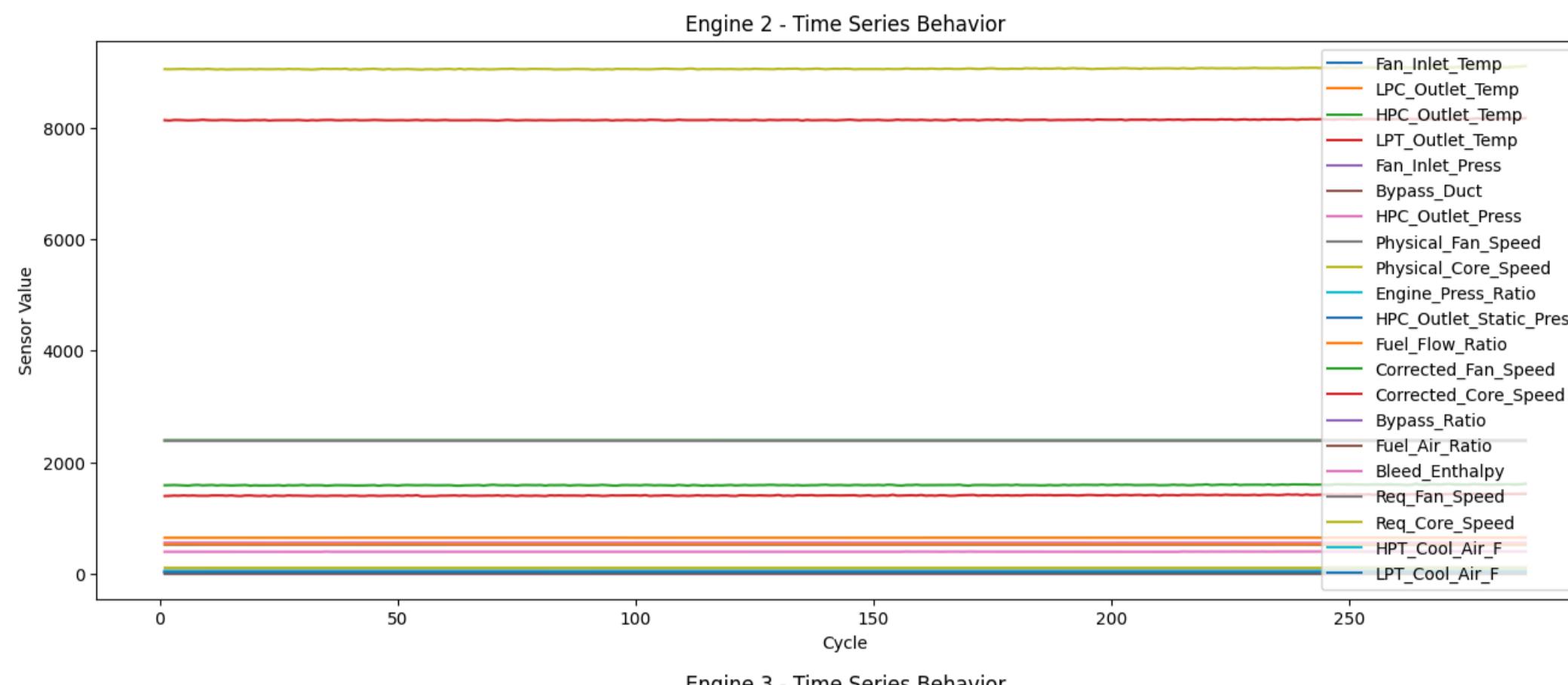
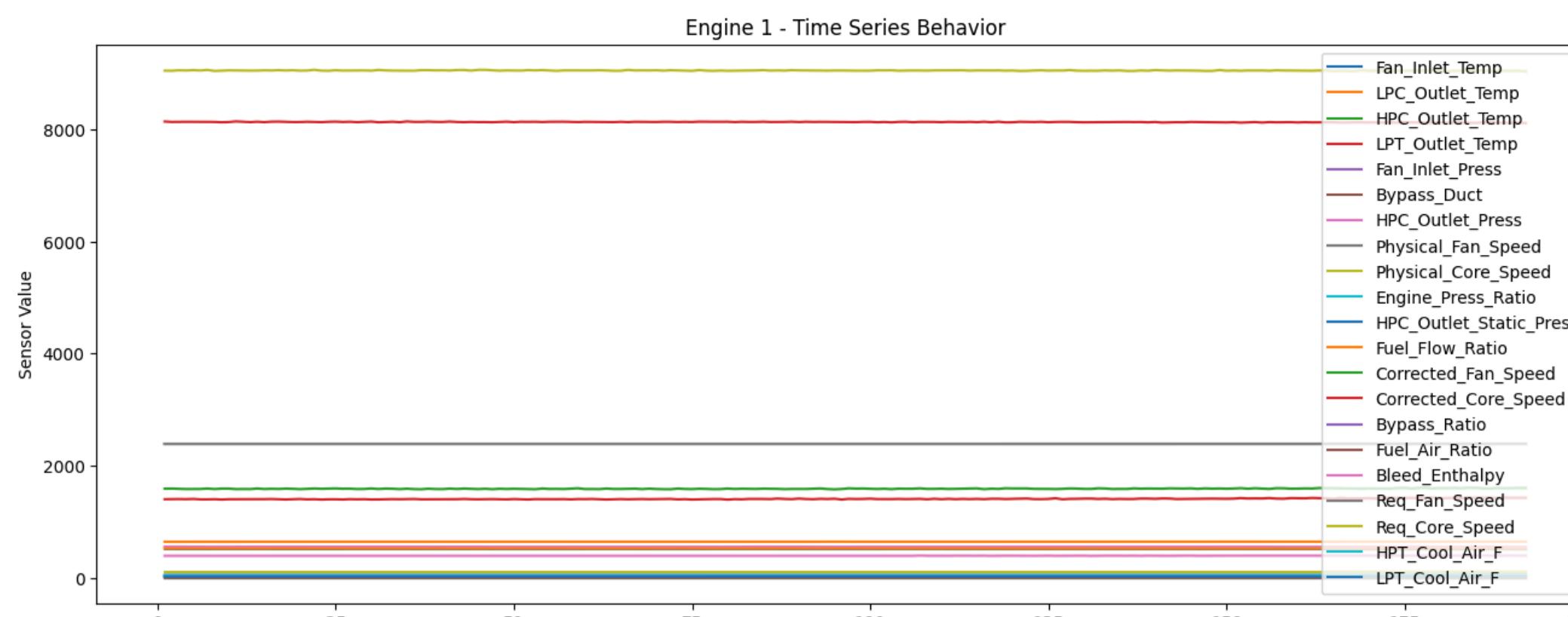


#### Time Series Behavior Analysis

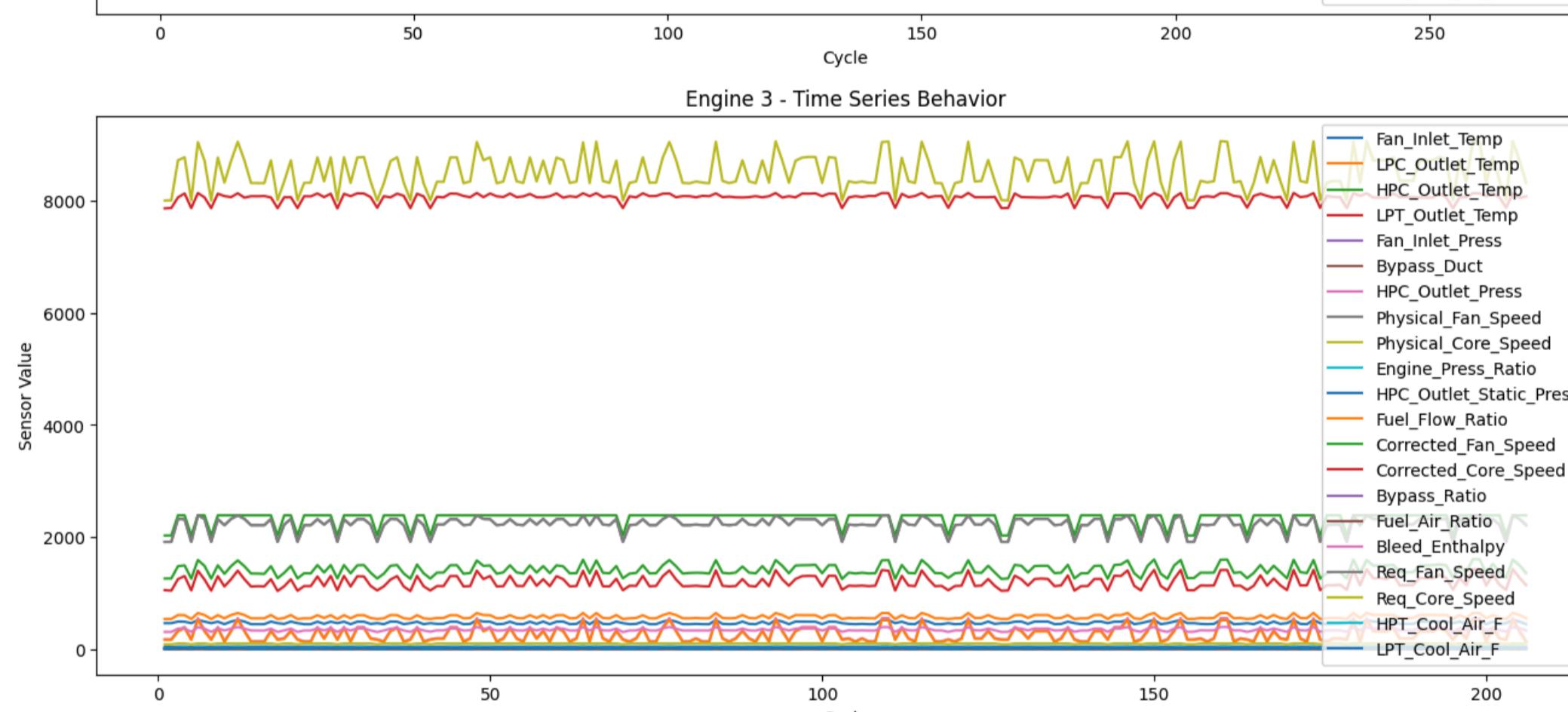
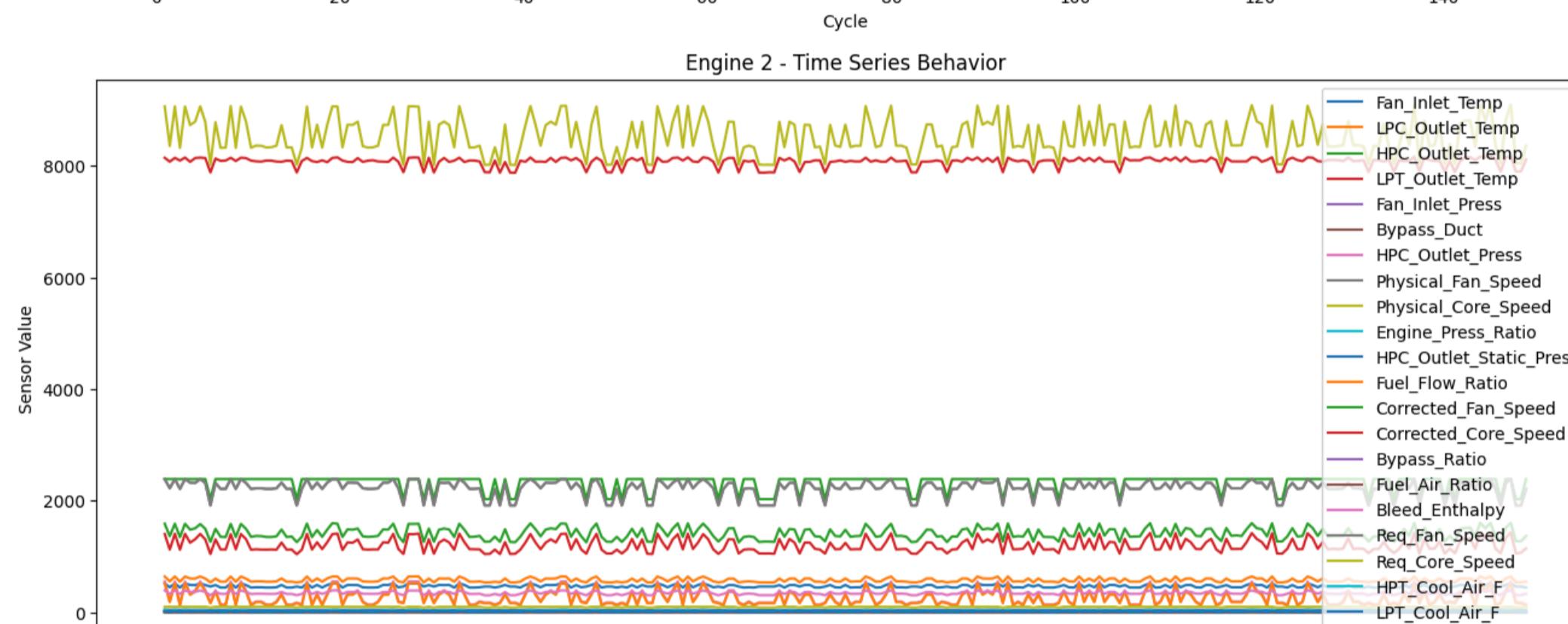
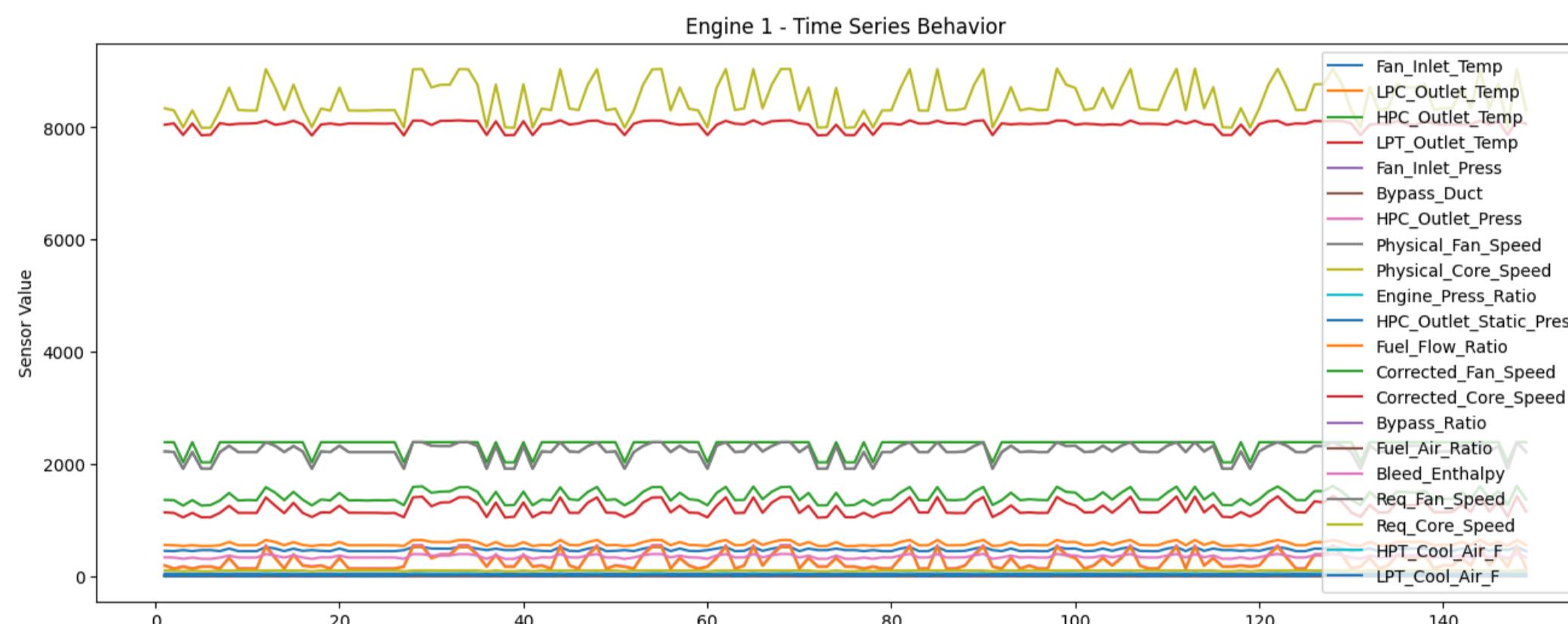
```
# 34/85
#train 1

sensors = [
    'Fan_Inlet_Temp',
    'LPC_Outlet_Temp',
    'HPC_Outlet_Temp',
    'LPT_Outlet_Temp',
    'Fan_Inlet_Press', # 5-9 (s1-s5)
    'Bypass_Duct',
    'HPC_Outlet_Press',
    'Physical_Fan_Speed',
    'Physical_Core_Speed',
    'Engine_Press_Ratio',
    'HPC_Outlet_Static_Press',
    'Fuel_Flow_Ratio',
    'Corrected_Fan_Speed',
    'Corrected_Core_Speed',
    'Bypass_Ratio', # 15-19 (s11-s15)
    'Fuel_Air_Ratio',
    'Bleed_Enthalpy',
    'Req_Fan_Speed',
    'Req_Core_Speed',
    'IPT_Cool_Air_F',
    'LPT_Cool_Air_F', # 20-25 (s16-s21, approx)
]

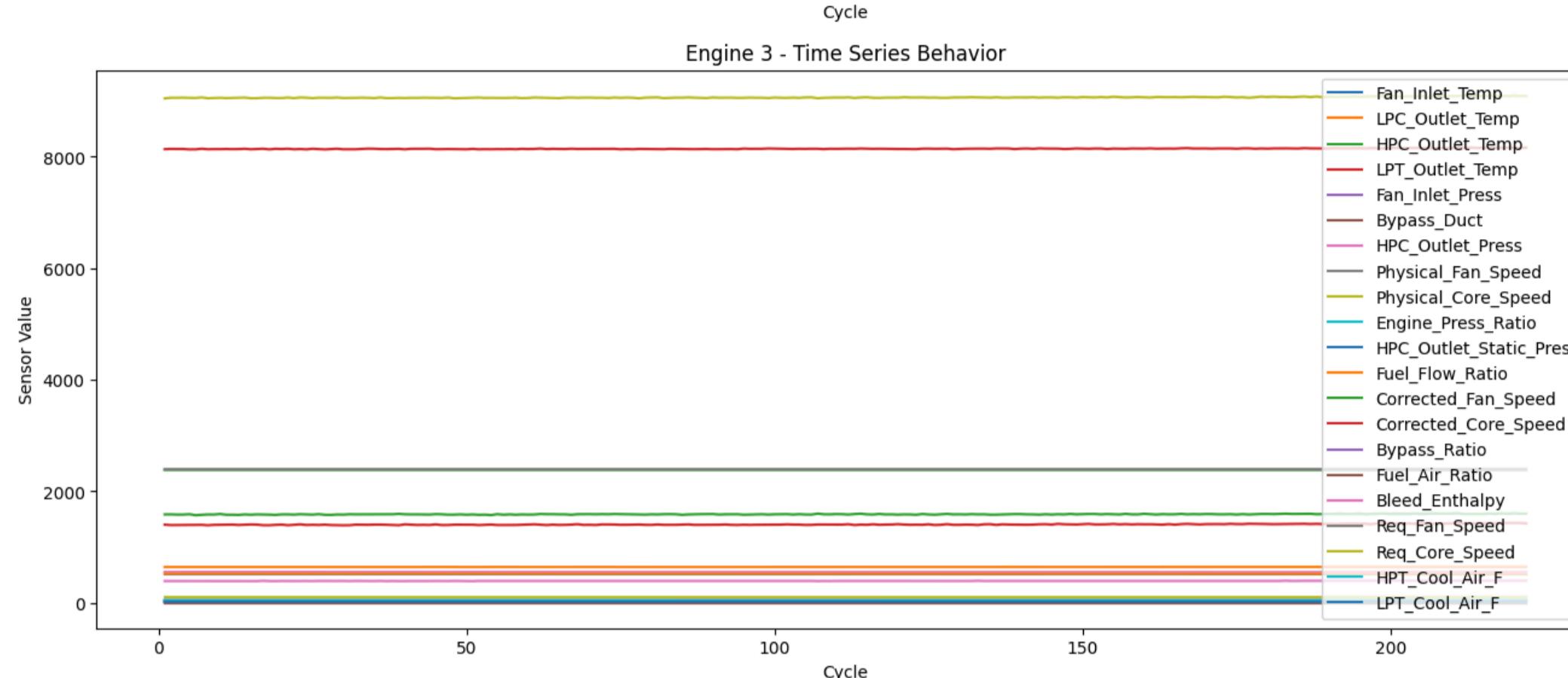
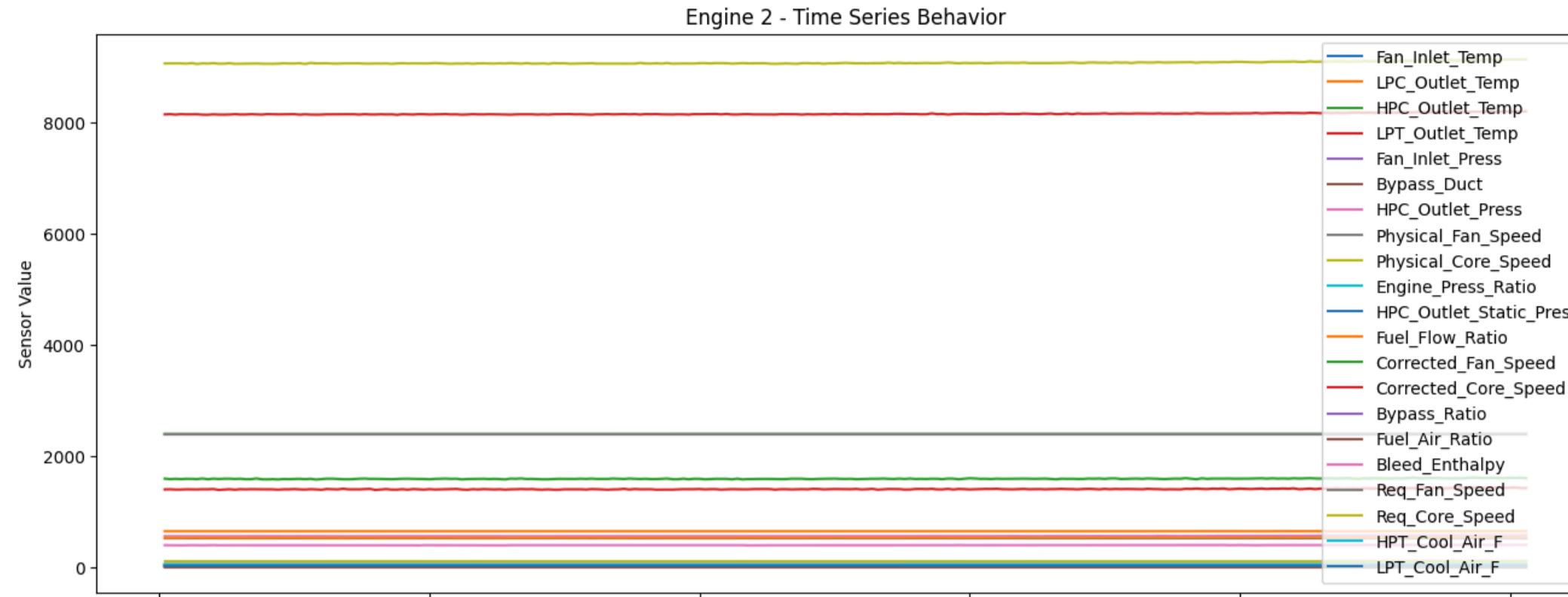
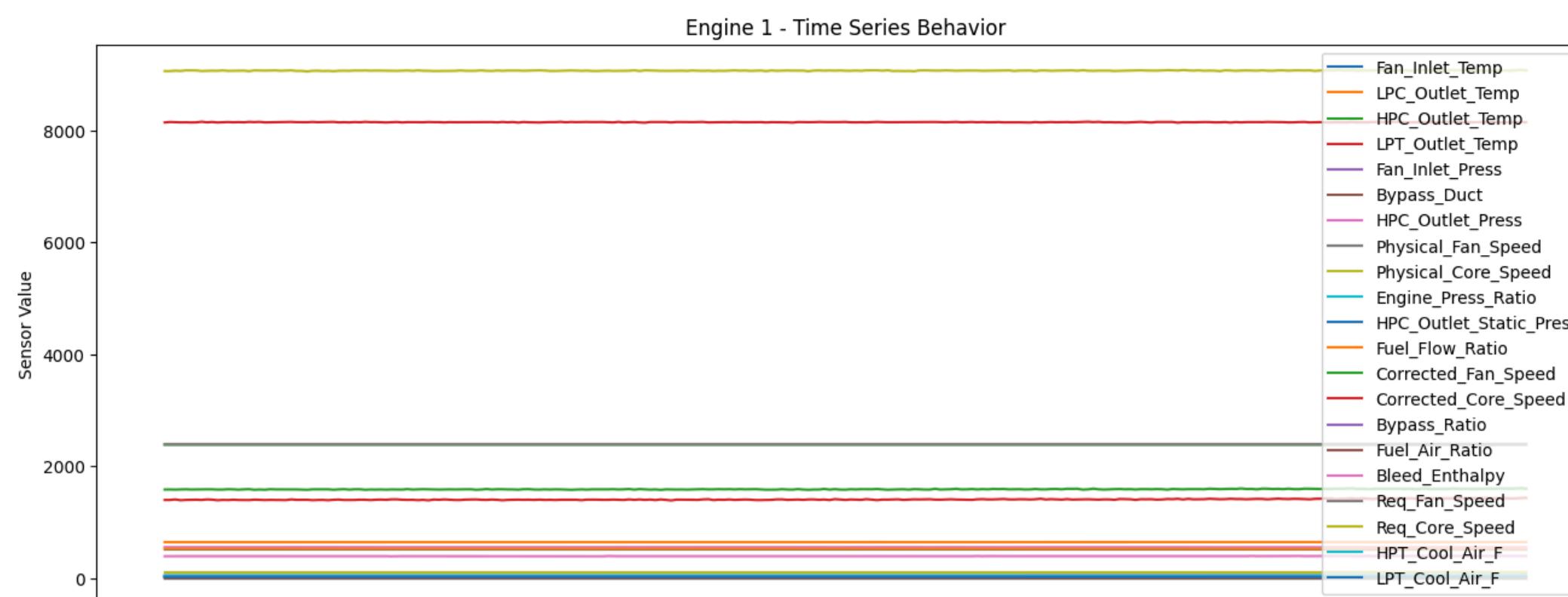
for engine_id in train1["unit"].unique():
    if engine_id in train1["unit"] == engine_id:
        plt.figure(figsize=(16, 6))
        for sensor in sensors:
            plt.plot(df[engine_id][sensor], df[engine_id][label=sensor])
        plt.title(f"Engine {engine_id} - Time Series Behavior")
        plt.xlabel("Time Step")
        plt.ylabel("Sensor Value")
        plt.legend()
        plt.show()
```



```
# 35/85
#train 2
for engine_id in train["unit"].unique()[:3]:
    df_engine = train[train["unit"] == engine_id]
    plt.figure(figsize=(16, 6))
    for sensor in sensors:
        plt.plot(df_engine["cycle"], df_engine[sensor], label=sensor)
    plt.title(f"Engine {engine_id} - Time Series Behavior")
    plt.xlabel("cycle")
    plt.ylabel("Sensor Value")
    plt.legend()
    plt.show()
```



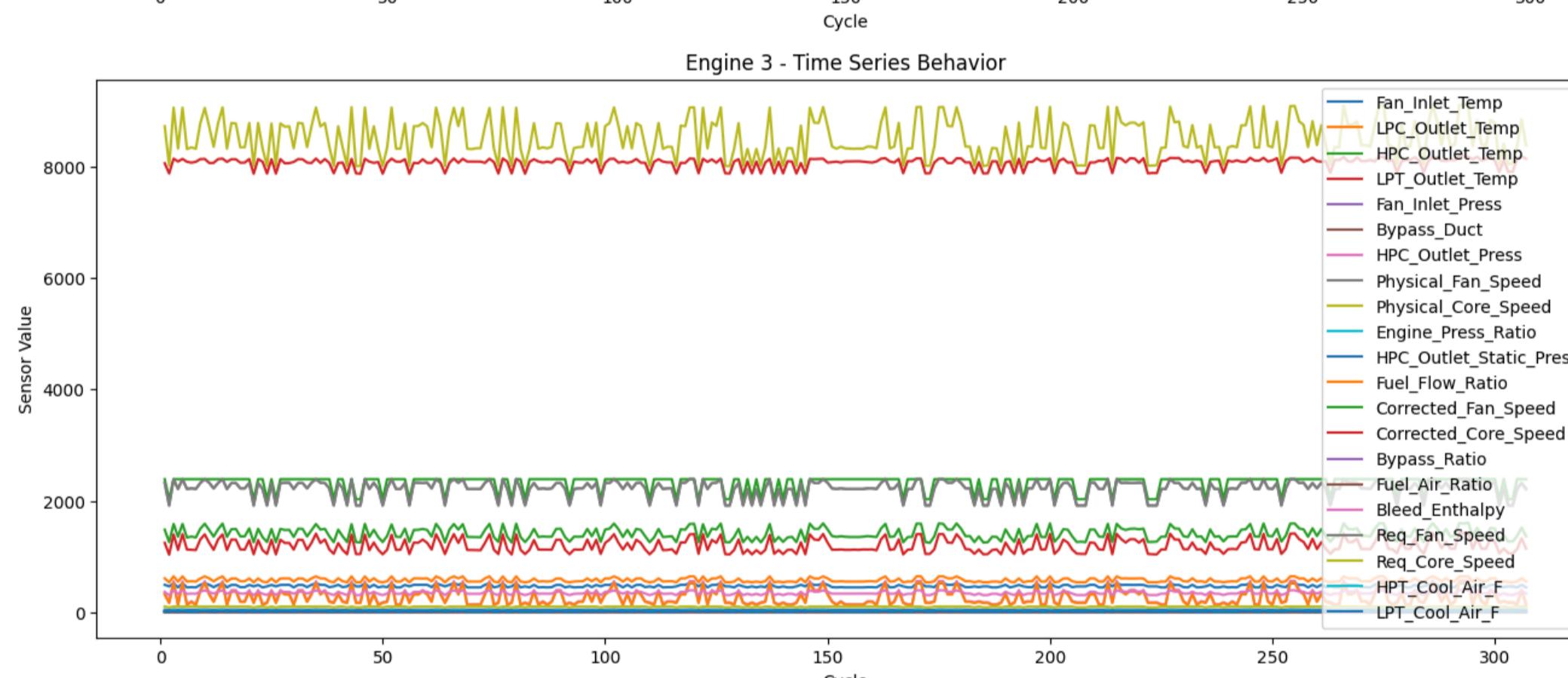
```
# 36/85
#train 3
for engine_id in train["unit"].unique()[:3]:
    df_engine = train[train["unit"] == engine_id]
    plt.figure(figsize=(16, 6))
    for sensor in sensors:
        plt.plot(df_engine["cycle"], df_engine[sensor], label=sensor)
    plt.title(f"Engine {engine_id} - Time Series Behavior")
    plt.xlabel("cycle")
    plt.ylabel("Sensor Value")
    plt.legend()
    plt.show()
```



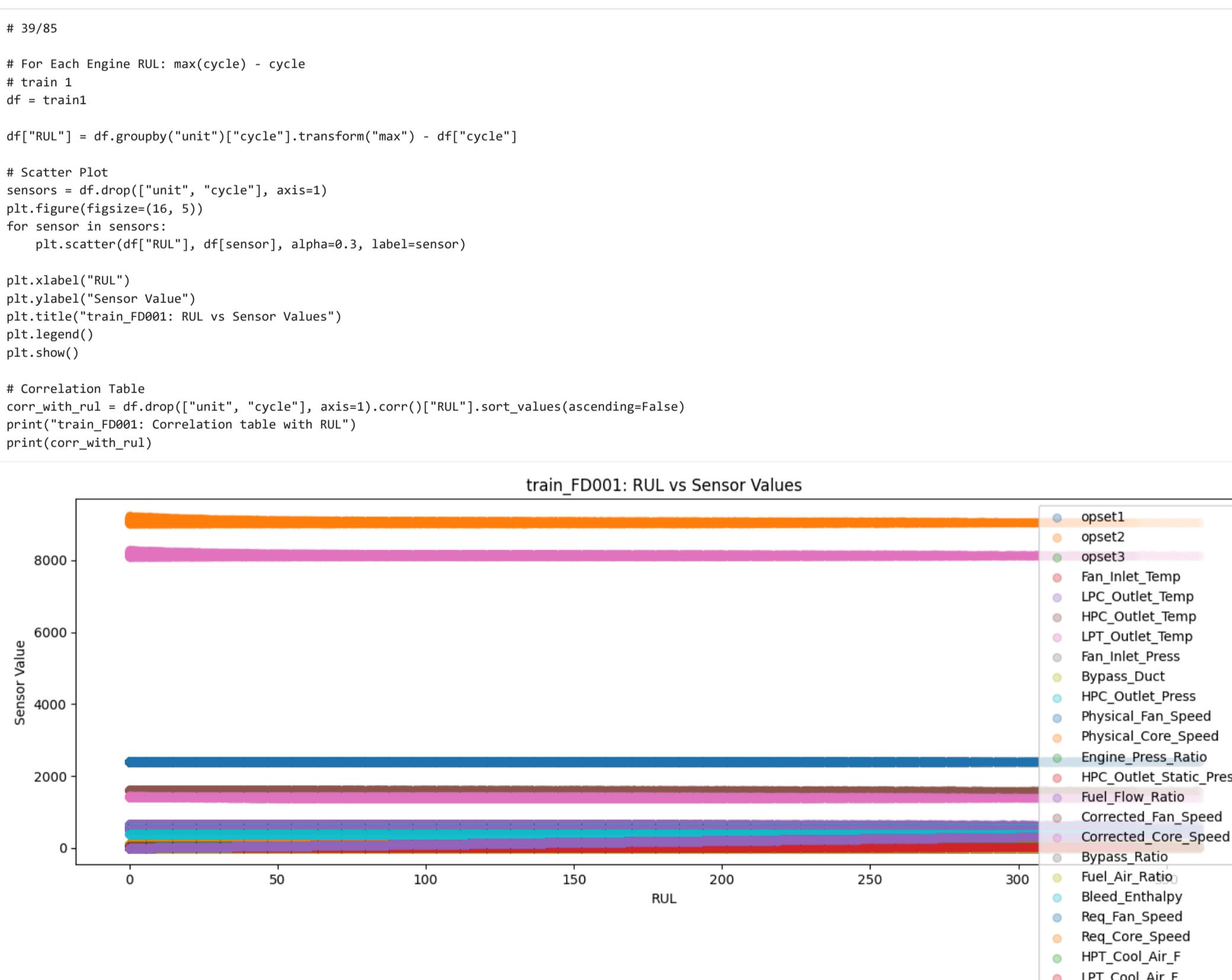
```
# 37/85
#train 4

for engine_id in train["unit"].unique()[:3]:
    df_engine = train[train["unit"] == engine_id]
    plt.figure(figsize=(16, 6))
    for sensor in sensors:
        plt.plot(df_engine["cycle"], df_engine[sensor], label=sensor)

    plt.title(f"Engine {engine_id} - Time Series Behavior")
    plt.xlabel("Cycle")
    plt.ylabel("Sensor Value")
    plt.legend()
    plt.show()
```



#### Relationship Analysis with Target (RUL)



train\_FD001: Correlation table with RUL

```
RUL      1.000000
Fuel_Flow_Ratio      0.863198
HPC_Outlet_Press      0.857223
LPT_Cool_Air_F      0.855662
oset1      0.853263
oset2      0.851562
oset3      0.849861
Fan_Inlet_Temp      0.848160
LPC_Outlet_Temp      0.846459
LPT_Outlet_Temp      0.844758
Fan_Inlet_Press      0.843057
Bypass_Duct      0.841356
HPC_Outlet_Press      0.839655
Physical_Fan_Speed      0.837954
Physical_Core_Speed      0.836253
Fuel_Flow_Ratio      0.834552
Corrected_Fan_Speed      0.832851
Corrected_Core_Speed      0.831150
Bypass_Ratio      0.829449
Fuel_Air_Ratio      0.827748
Bleed_Enthalpy      0.826047
Req_Fan_Speed      0.824346
Req_Core_Speed      0.822645
HPT_Cool_Air_F      0.821944
LPT_Cool_Air_F      0.820243
RUL      0.819543
```

```
# 40/85
# train 2

df["RUL"] = df.groupby("unit")["cycle"].transform("max") - df["cycle"]

df["RUL"] = df.groupby("unit")["cycle"].transform("max") - df["cycle"]

# Scatter Plot
sensors = df.drop(["unit", "cycle"], axis=1)
plt.figure(figsize=(16, 5))
for sensor in sensors:
    plt.scatter(df["RUL"], df[sensor], alpha=0.3, label=sensor)

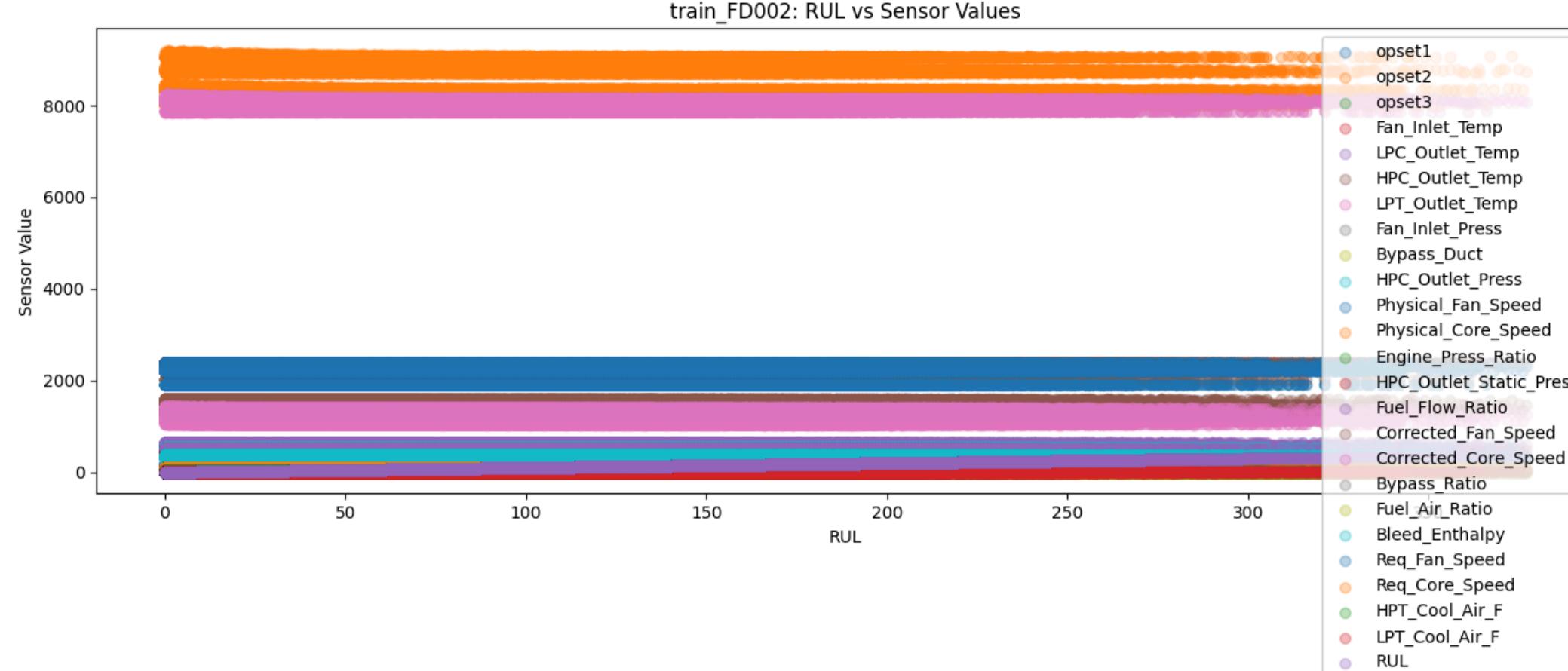
plt.xlabel("RUL")
plt.ylabel("Sensor Value")
plt.title("train_FD001: RUL vs Sensor Values")
plt.show()
```

```

plt.xlabel("RUL")
plt.ylabel("Sensor Value")
plt.title("train_FD002: RUL vs Sensor Values")
plt.legend()
plt.show()

# Correlation Table
cor_with_rul = df.drop(["unit", "cycle"], axis=1).corr()["RUL"].sort_values(ascending=False)
print("train_FD002: Correlation table with RUL")
print(cor_with_rul)

```



```

train_FD002: Correlation table with RUL
RUL           1.000000
HPT_Cool_Air_F      0.000165
LPT_Cool_Air_F      0.000165
Req_Core_Speed      0.0005761
opset1            0.0005761
opset2            0.0005245
opset3            0.0004785
Req_Fan_Speed      0.0004785
Physical_Fan_Speed 0.0004386
Engine_Press_Ratio 0.0004386
HPC_Outlet_Press    0.00042486
Fuel_Air_Ratio      0.00042302
opset1            0.0001135
opset2            0.000047
opset3            0.000047
Fan_Inlet_Temp     -0.0004956
Bypass_Duct        -0.0004956
Fan_Inlet_Press    -0.000758
LPC_Outlet_Temp    -0.000758
Physical_Core_Speed -0.015380
HPC_Outlet_Temp    -0.015942
Bypass_Ratio       -0.015942
Bypass_Duct        -0.0184555
LPT_Outlet_Temp    -0.0184555
Corrected_Fan_Speed 0.0184555
Corrected_Core_Speed 0.0184555
HPC_Outlet_Static_Press -0.046952
Fuel_Air_Ratio      -0.071352
Name: RUL, dtype: float64

```

```

# 4/1/85
# train 3
df = train3

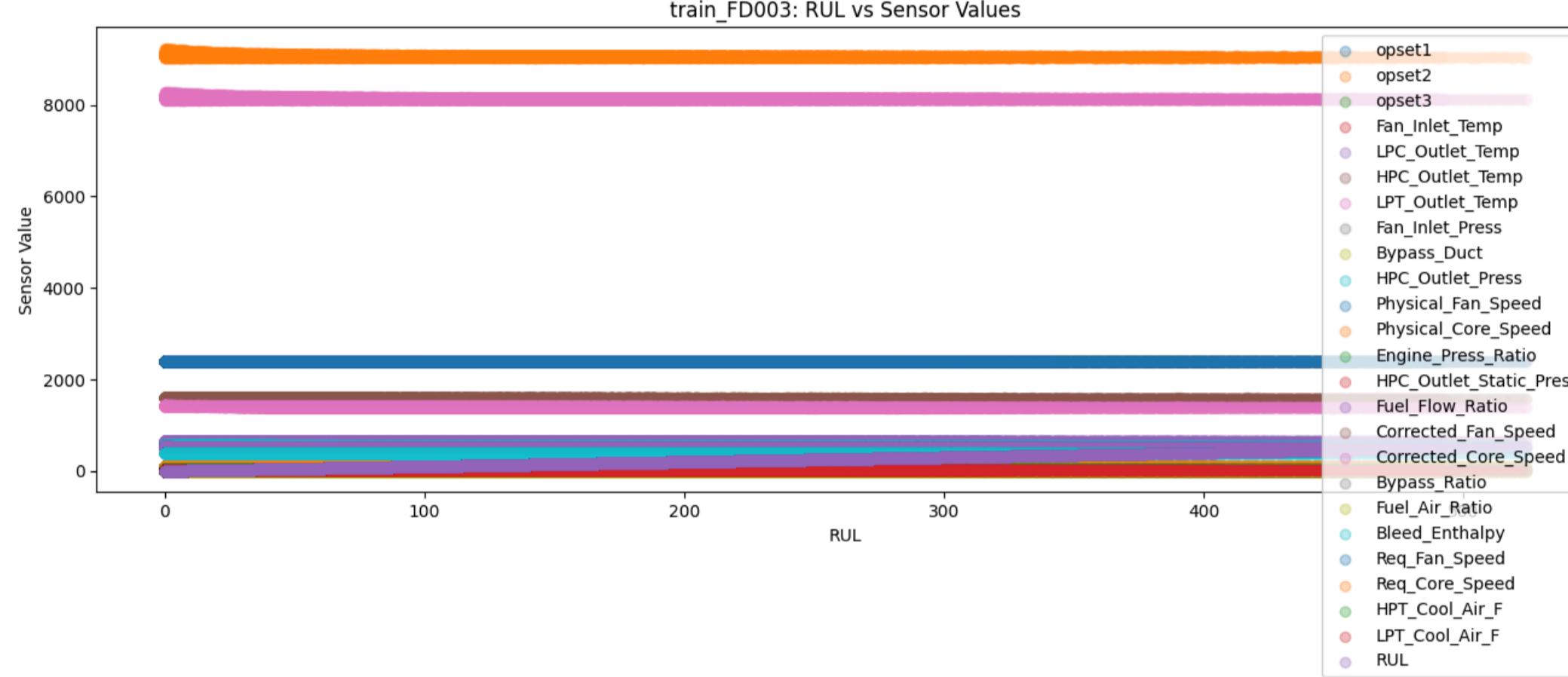
df[["RUL"]] = df.groupby("unit")["cycle"].transform("max") - df[["cycle"]]

# Scatter Plot
sensors = df.drop(["unit", "cycle"], axis=1)
plt.figure(figsize=(16, 5))
for sensor in sensors:
    plt.scatter(df["RUL"], df[sensor], alpha=0.3, label=sensor)

plt.xlabel("RUL")
plt.ylabel("Sensor Value")
plt.title("train_FD003: RUL vs Sensor Values")
plt.legend()
plt.show()

# Correlation Table
cor_with_rul = df.drop(["unit", "cycle"], axis=1).corr()["RUL"].sort_values(ascending=False)
print("train_FD003: Correlation table with RUL")
print(cor_with_rul)

```



```

train_FD003: Correlation table with RUL
RUL           1.000000
HPT_Cool_Air_F      0.017782
LPT_Cool_Air_F      0.013465
opset1            0.013465
opset2            -0.001284
opset3            -0.001284
Bypass_Ratio       -0.016581
Bypass_Duct        -0.016581
HPC_Outlet_Press    0.015948
Fuel_Flow_Ratio     -0.329452
Engine_Press_Ratio 0.0129452
opset1            0.491887
Corrected_Core_Speed 0.551843
Physical_Core_Speed 0.551843
LPC_Outlet_Temp    0.624236
HPC_Outlet_Temp    0.649236
Bleed_Enthalpy     -0.649236
Physical_Fan_Speed 0.656348
opset2            0.656348
Corrected_Fan_Speed 0.657224
LPT_Outlet_Temp    0.657224
HPC_Outlet_Static_Press -0.660224
opset3            NaN
Fan_Inlet_Temp     NaN
Fan_Inlet_Press    NaN
Fuel_Air_Ratio      NaN
Req_Fan_Speed      NaN
Req_Core_Speed      NaN
Name: RUL, dtype: float64

```

```

# 4/2/85
# train 4
df = train4

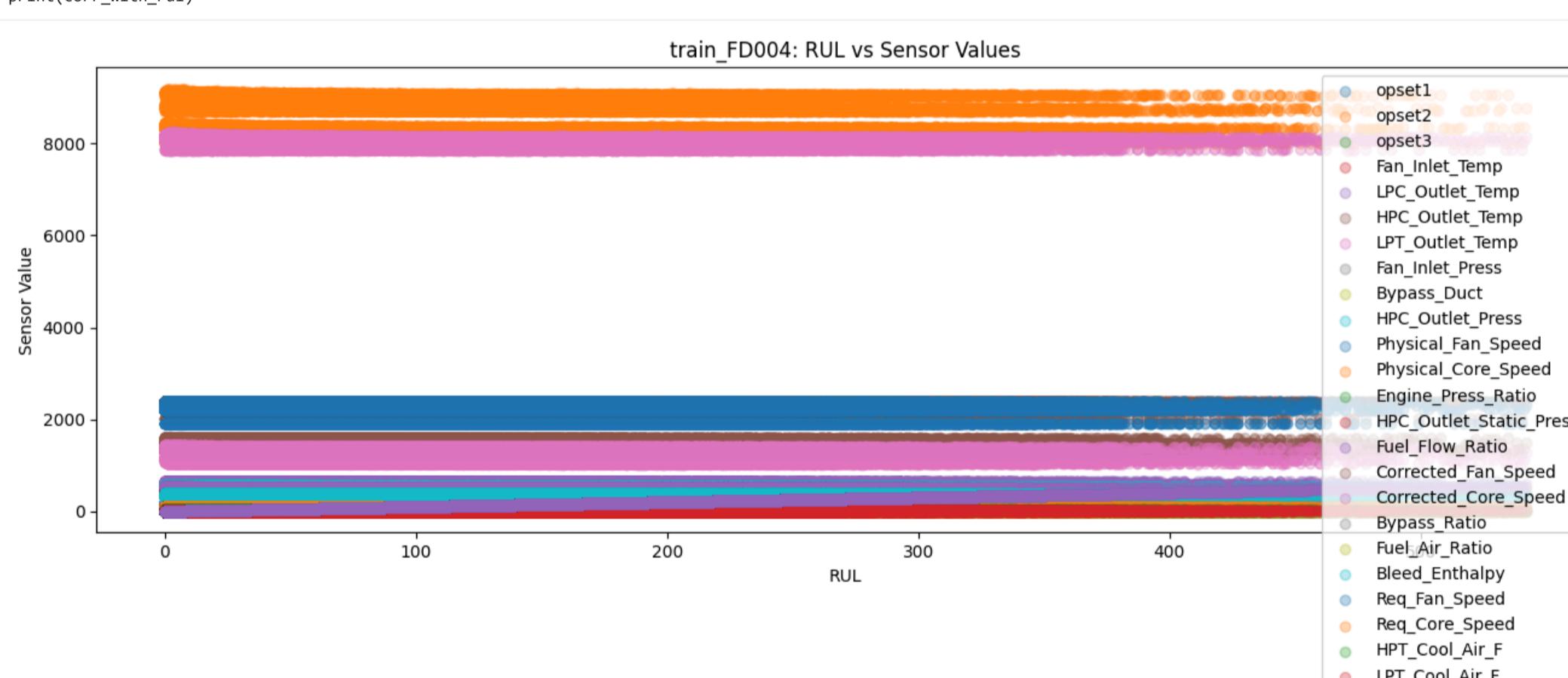
df[["RUL"]] = df.groupby("unit")["cycle"].transform("max") - df[["cycle"]]

# Scatter Plot
sensors = df.drop(["unit", "cycle"], axis=1)
plt.figure(figsize=(16, 5))
for sensor in sensors:
    plt.scatter(df["RUL"], df[sensor], alpha=0.3, label=sensor)

plt.xlabel("RUL")
plt.ylabel("Sensor Value")
plt.title("train_FD004: RUL vs Sensor Values")
plt.legend()
plt.show()

# Correlation Table
cor_with_rul = df.drop(["unit", "cycle"], axis=1).corr()["RUL"].sort_values(ascending=False)
print("train_FD004: Correlation table with RUL")
print(cor_with_rul)

```



```

train_FD004: Correlation table with RUL
RUL           1.000000
HPT_Cool_Air_F      0.002812
LPT_Cool_Air_F      0.000101
Req_Fan_Speed      0.000765
opset3            0.002383
opset2            0.000393
Physical_Fan_Speed 0.000393
Physical_Core_Speed 0.000393
Fan_Inlet_Temp     0.0003886
Fan_Inlet_Press    0.0003886
Corrected_Fan_Speed 0.0001561
Bypass_Duct        0.0001561
HPC_Outlet_Temp    0.0001349
Fuel_Flow_Ratio     -0.0001639
opset2            -0.0002286
opset1            -0.0002286
opset3            -0.0003957
LPC_Outlet_Temp    -0.0004443
Engine_Press_Ratio -0.0004443
Physical_Core_Speed -0.024727
HPC_Outlet_Temp    -0.019294
Bypass_Ratio       -0.019294
opset1            -0.0458811
LPT_Outlet_Temp    -0.0458811
Fuel_Air_Ratio      -0.053884
HPC_Outlet_Static_Press -0.053884
Corrected_Core_Speed -0.079126
Name: RUL, dtype: float64

```

#### Comparison of All Sensors Across All Datasets

```

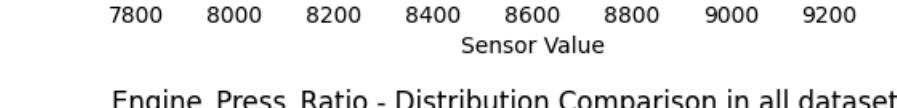
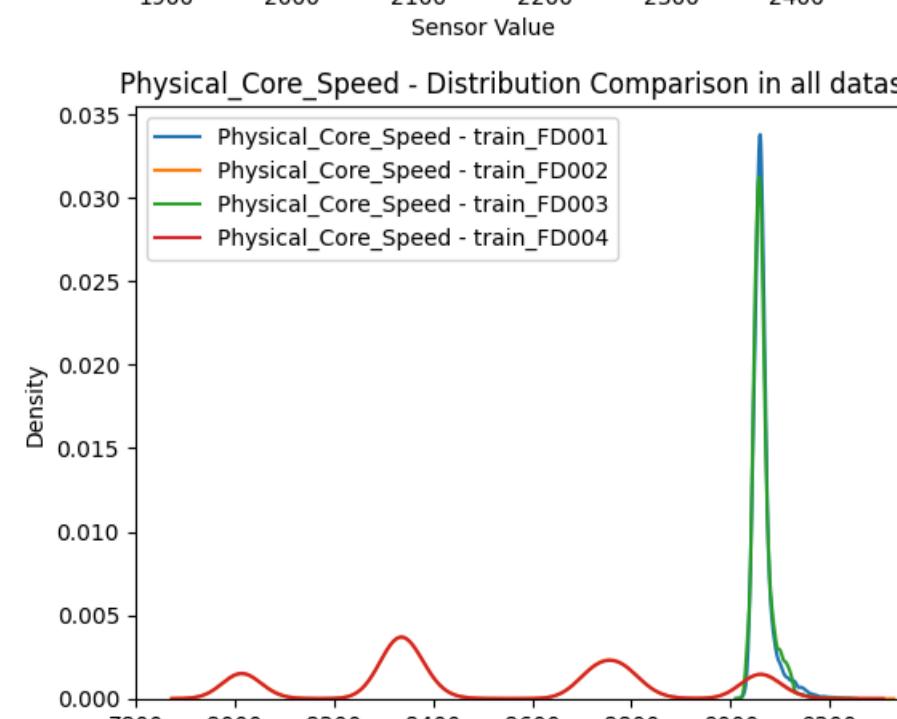
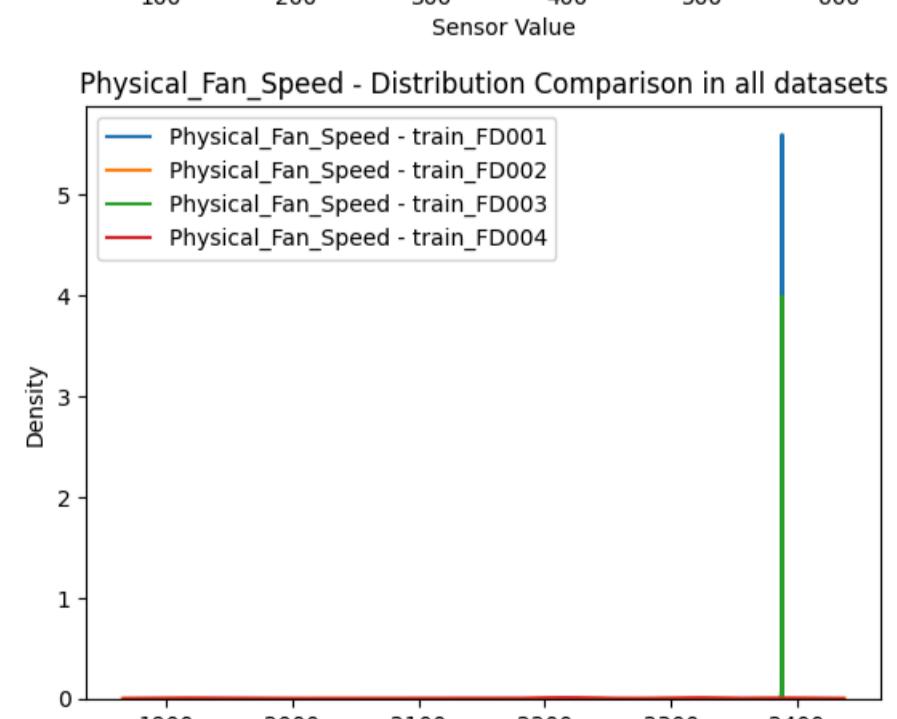
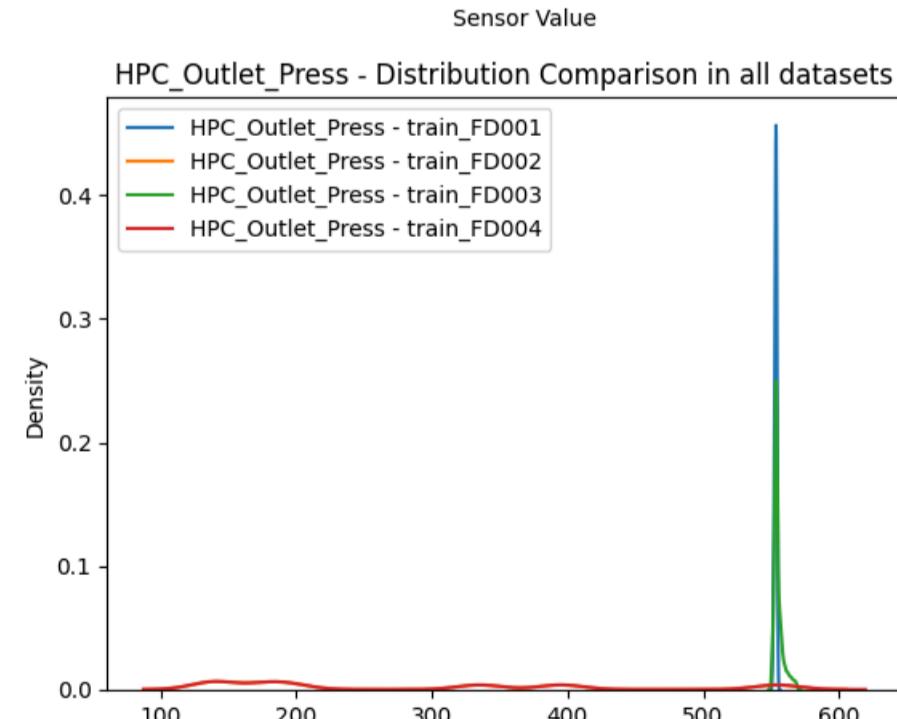
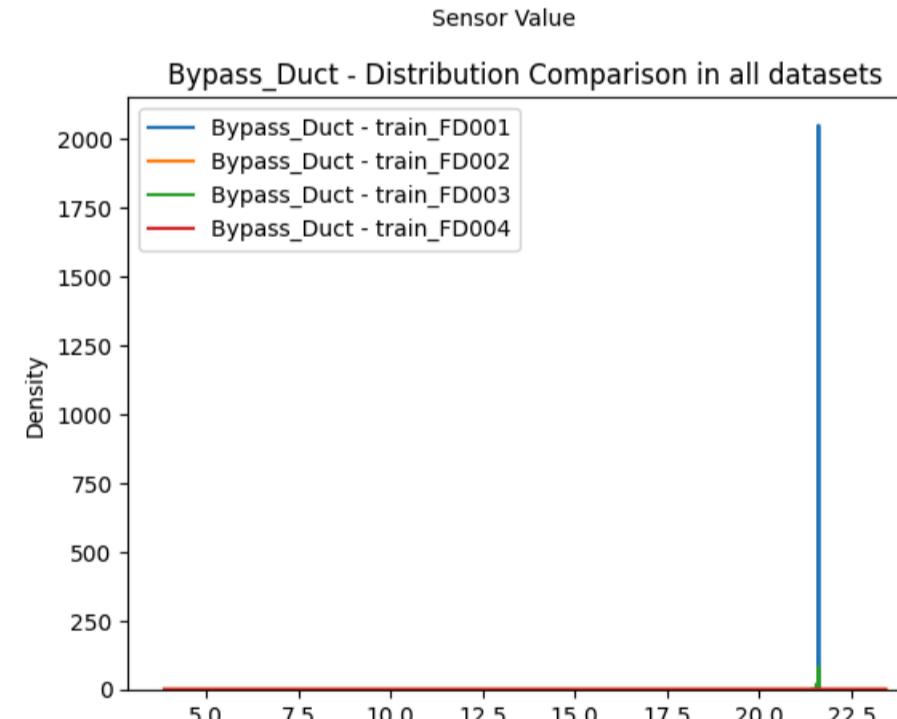
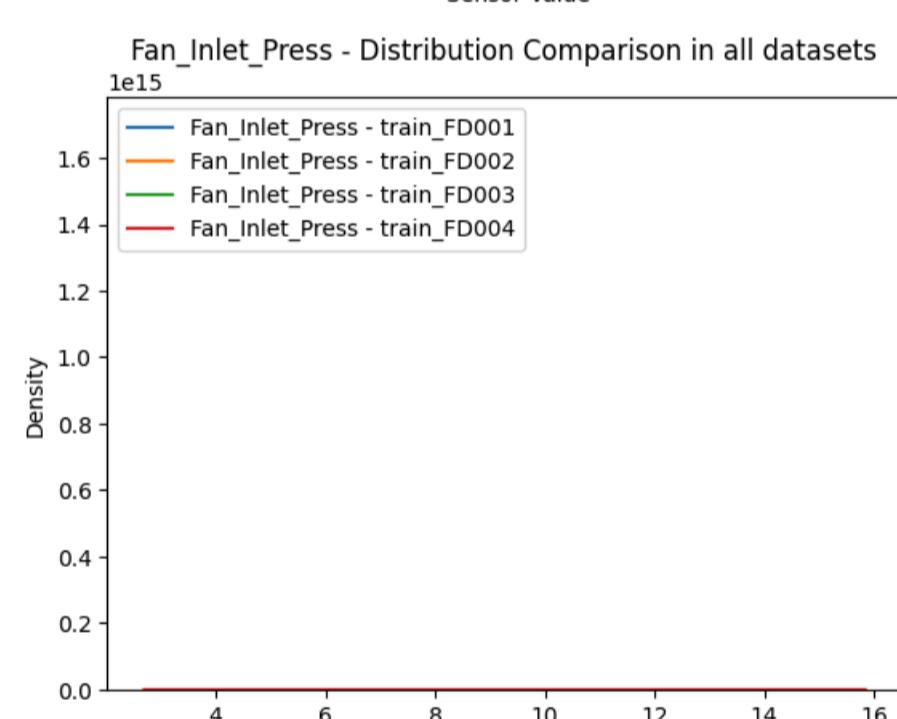
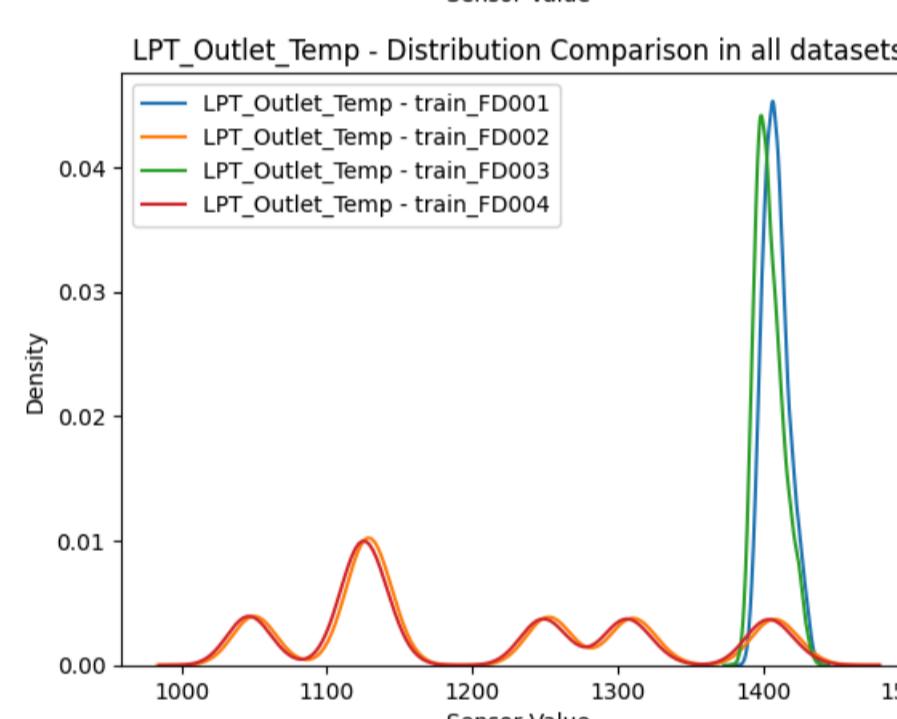
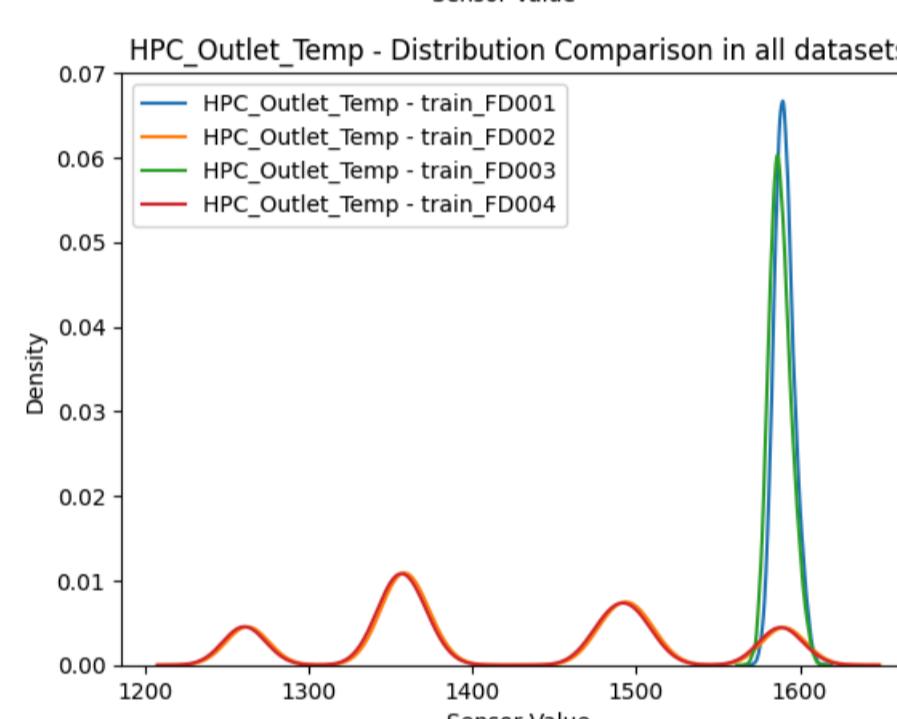
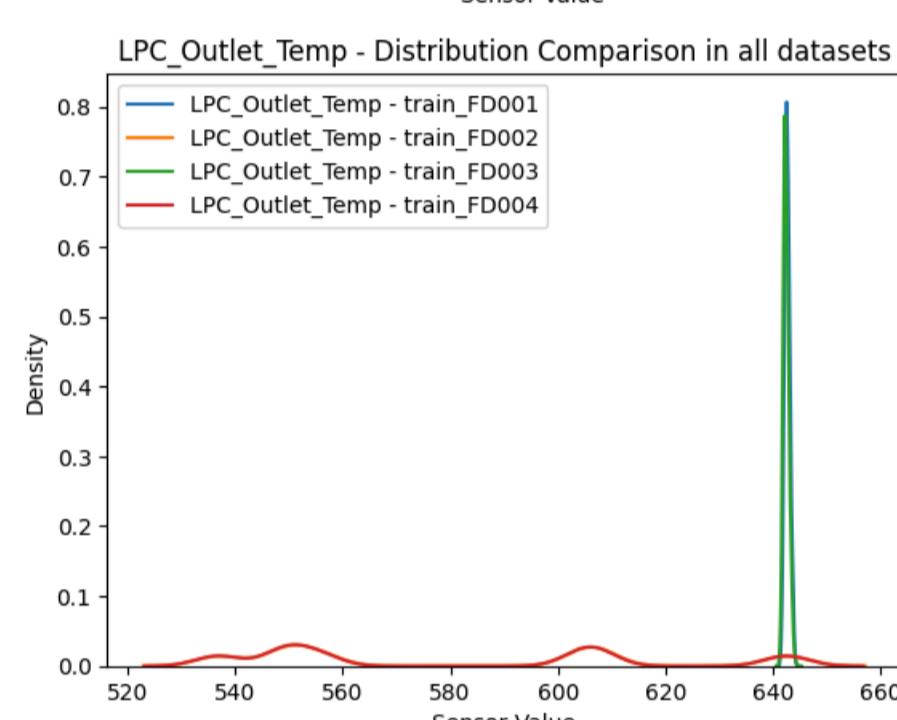
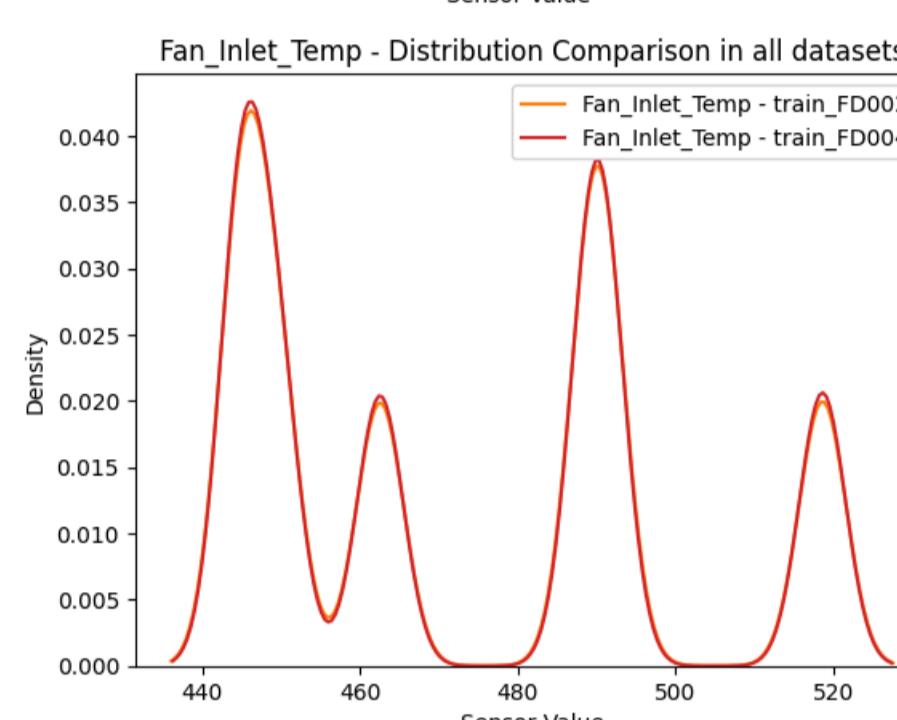
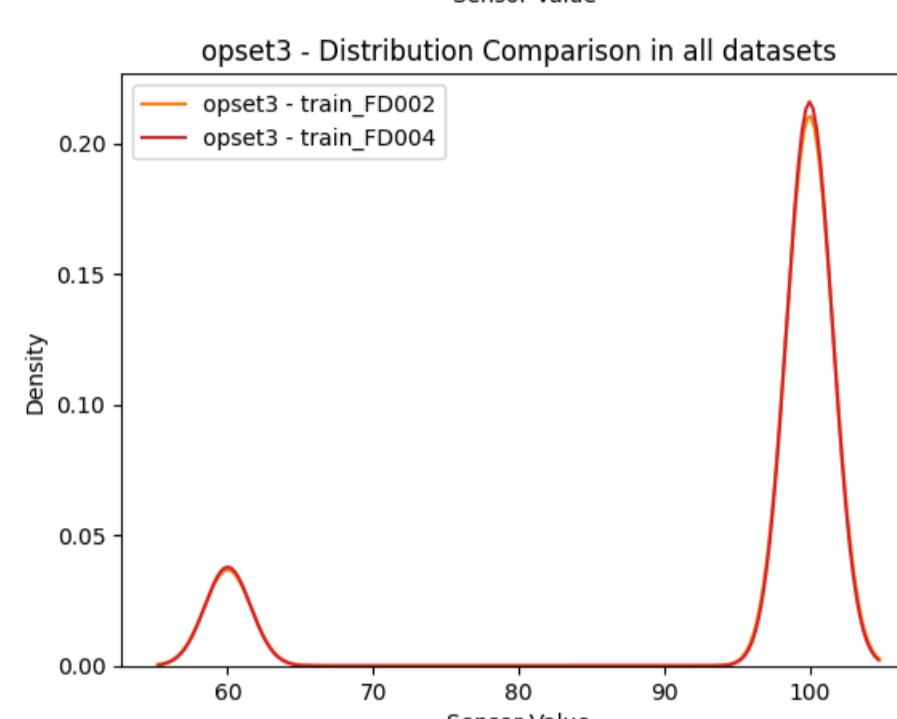
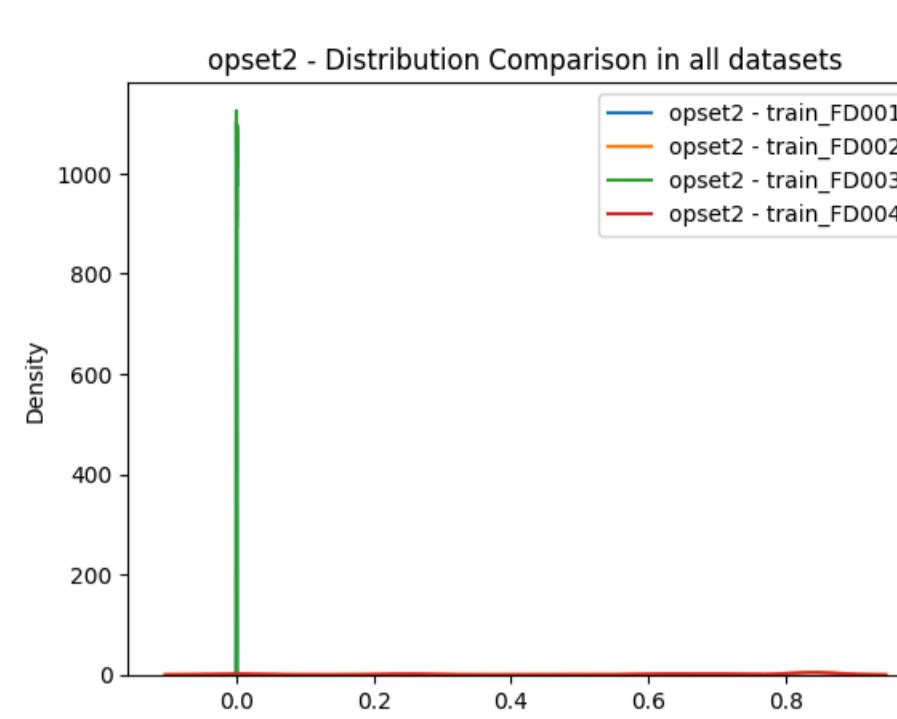
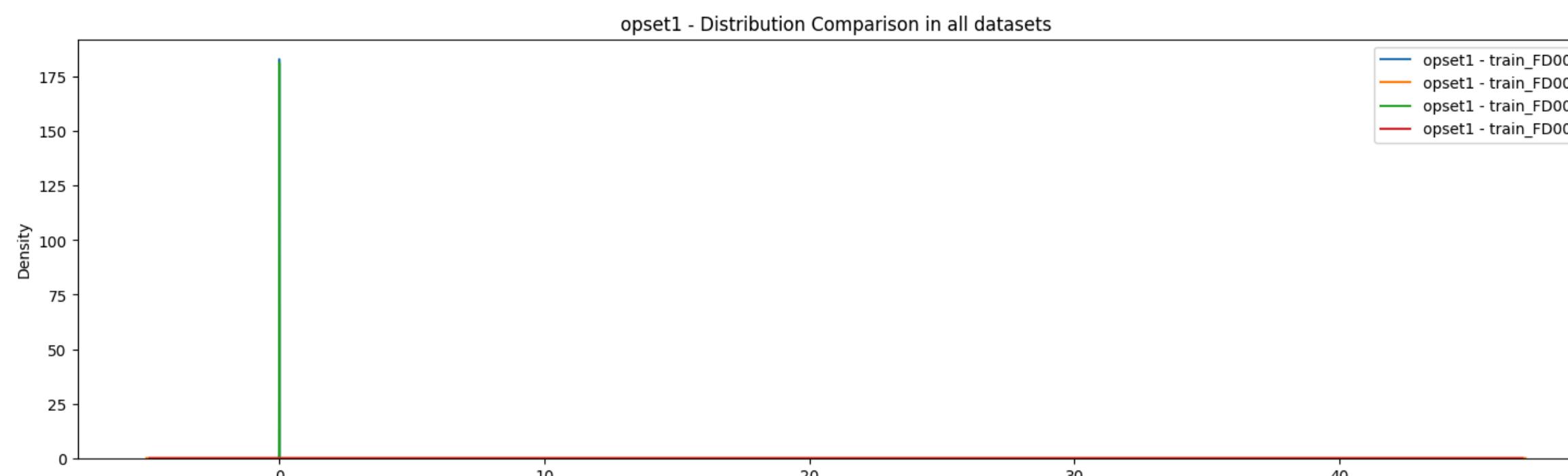
# 44/85
sensors = df.drop(["unit", "cycle"], axis=1)
plt.figure(figsize=(18, 5))

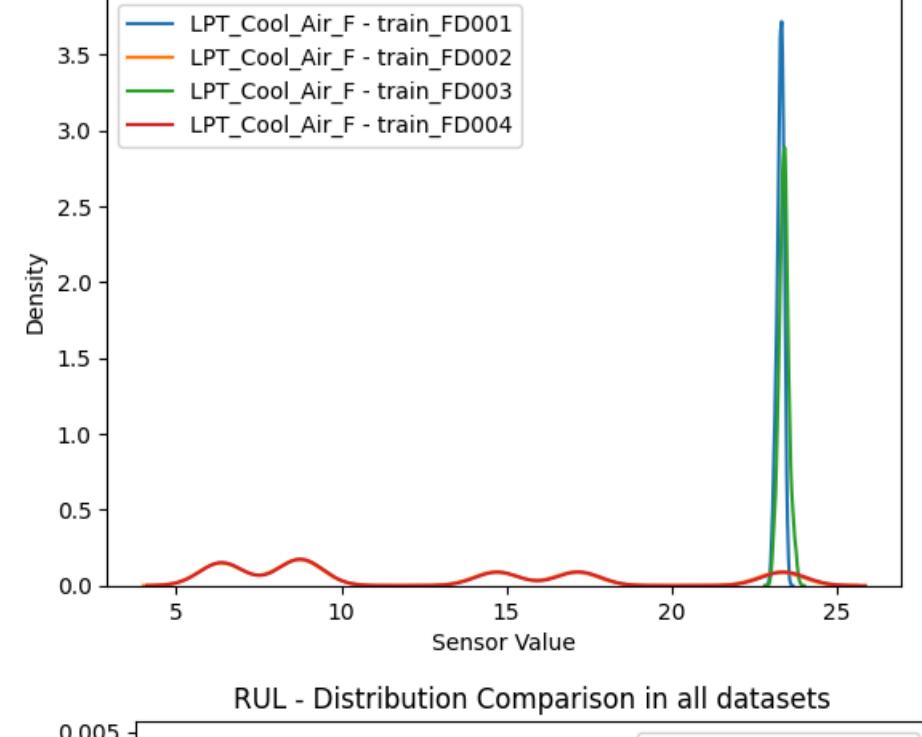
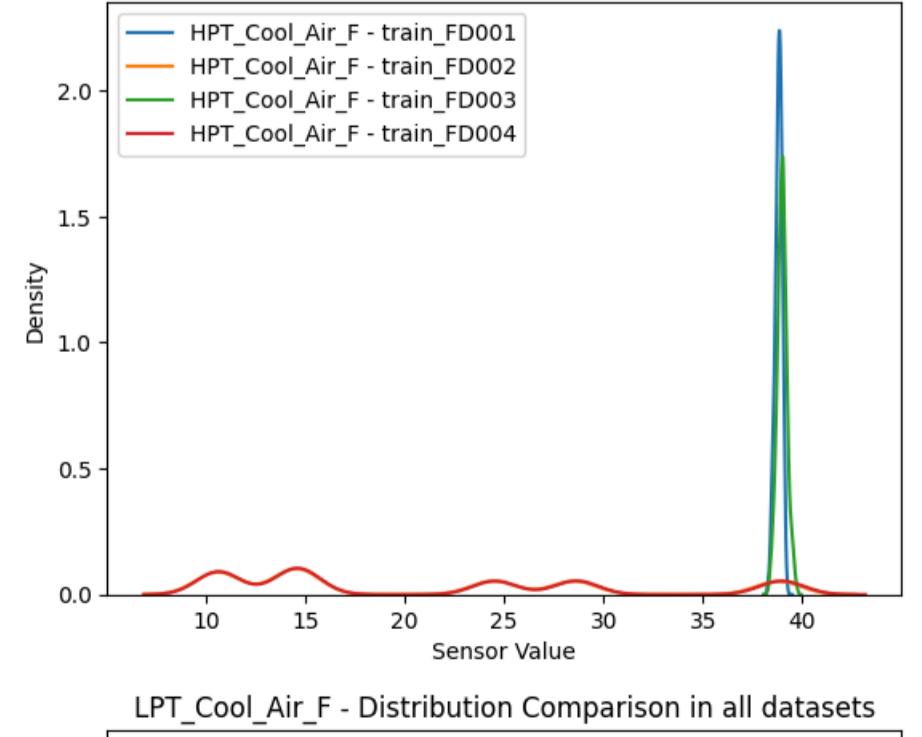
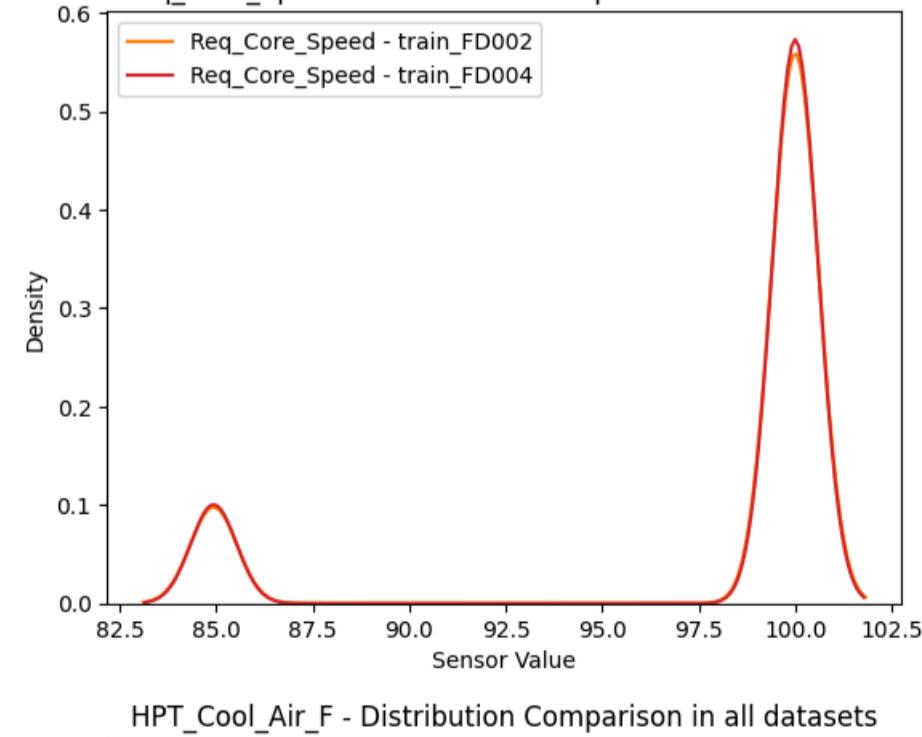
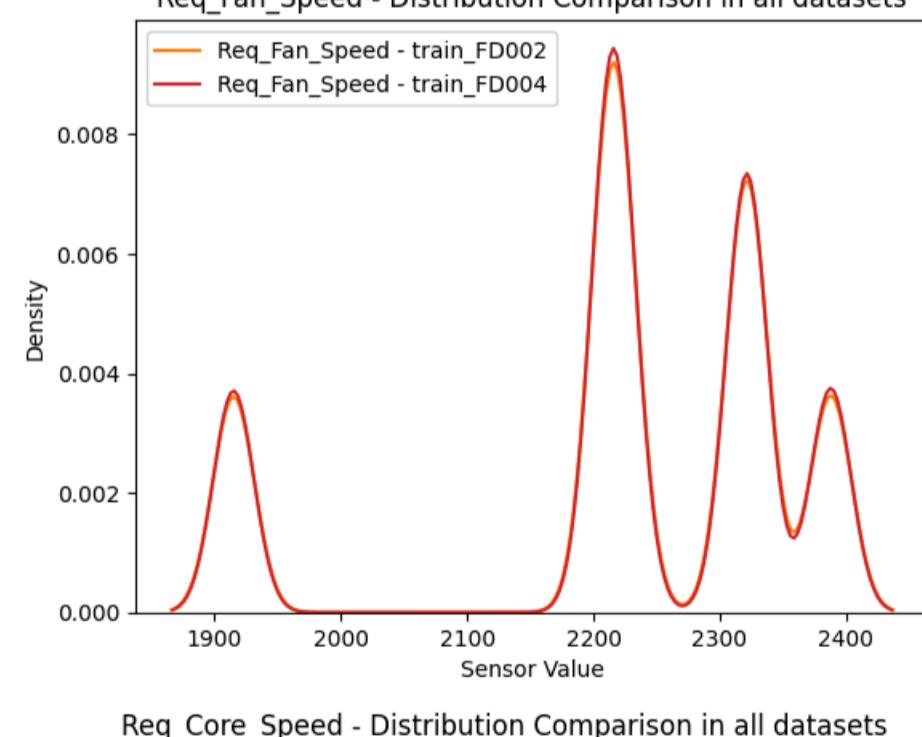
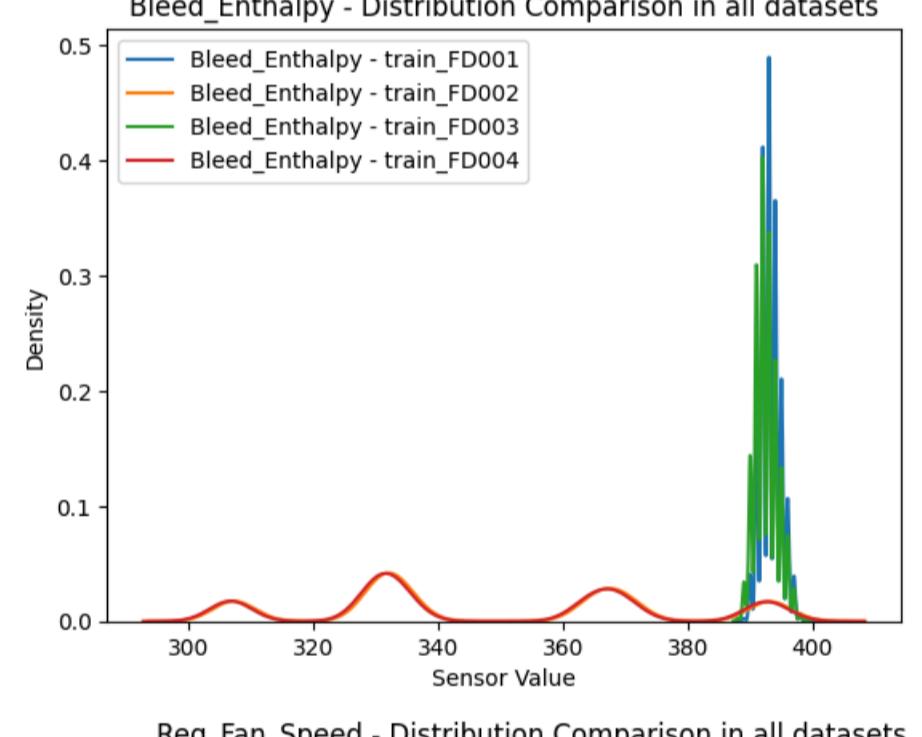
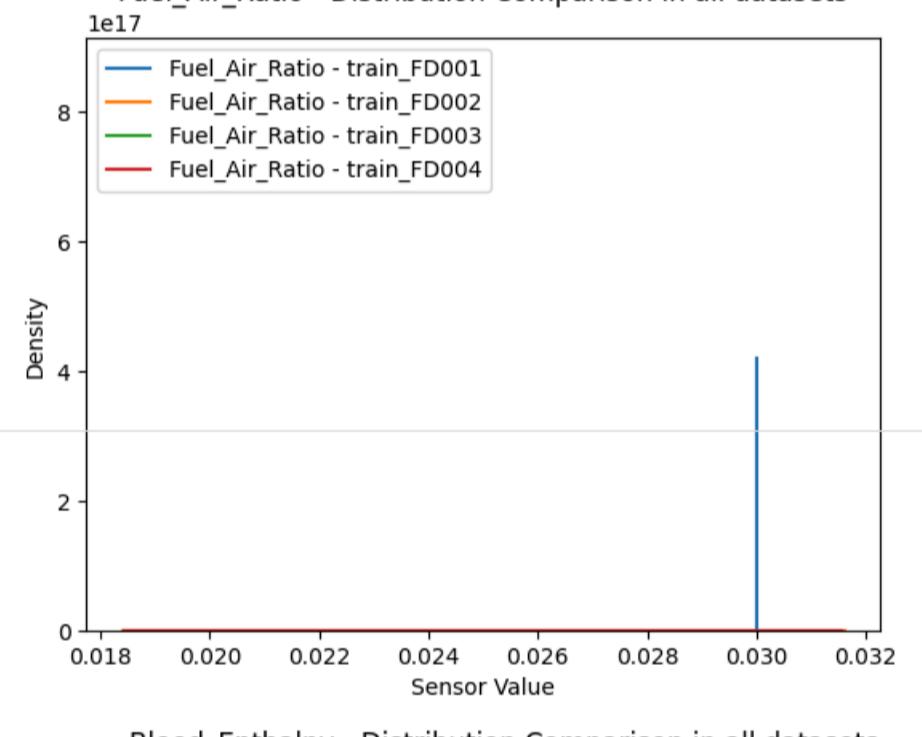
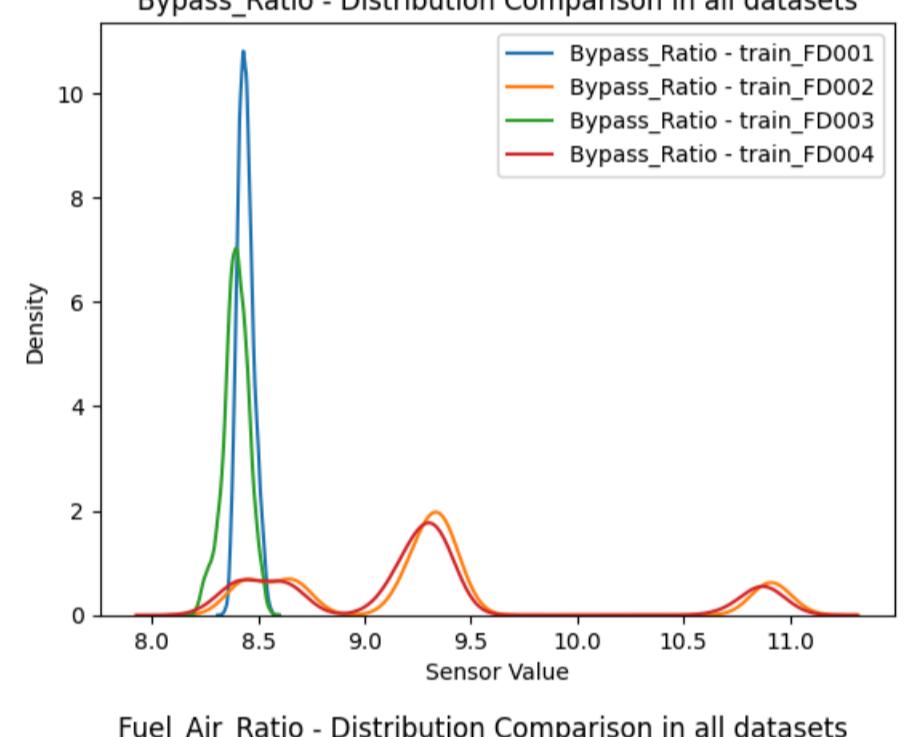
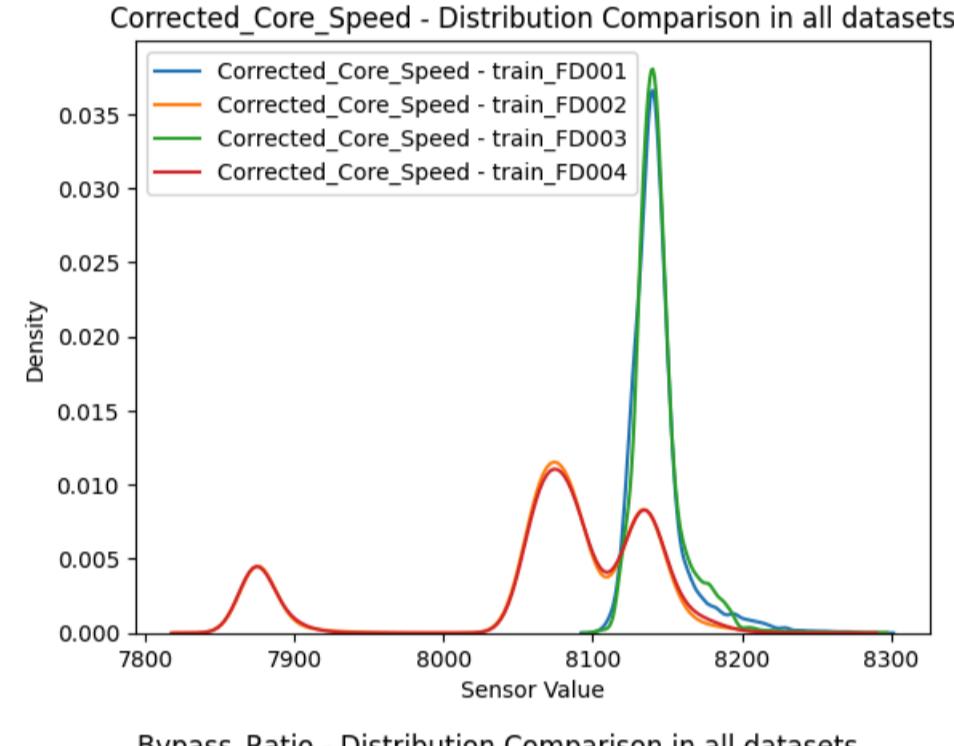
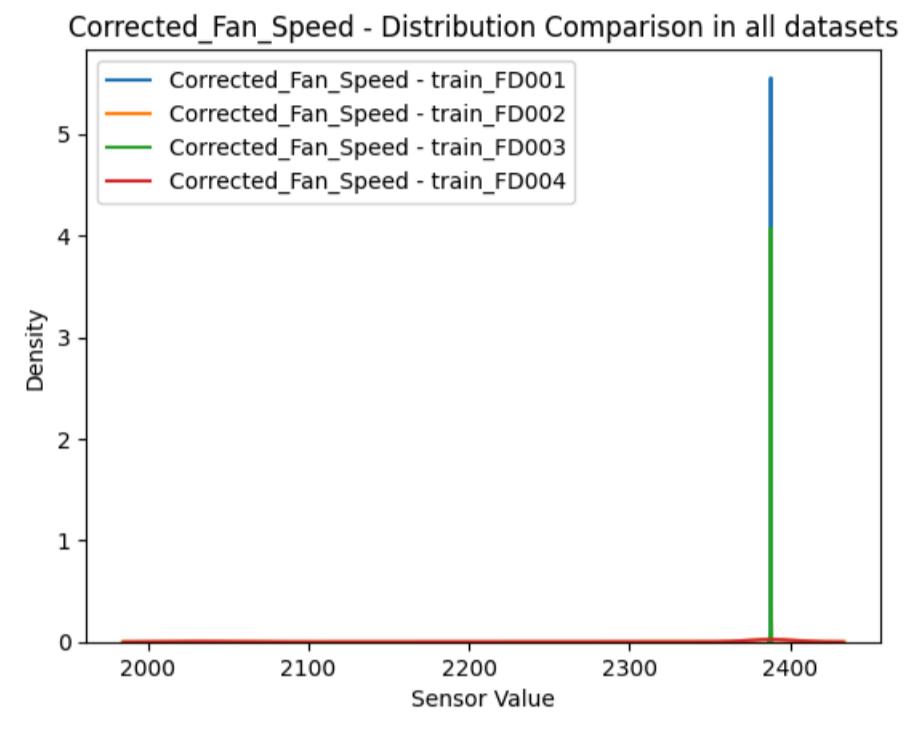
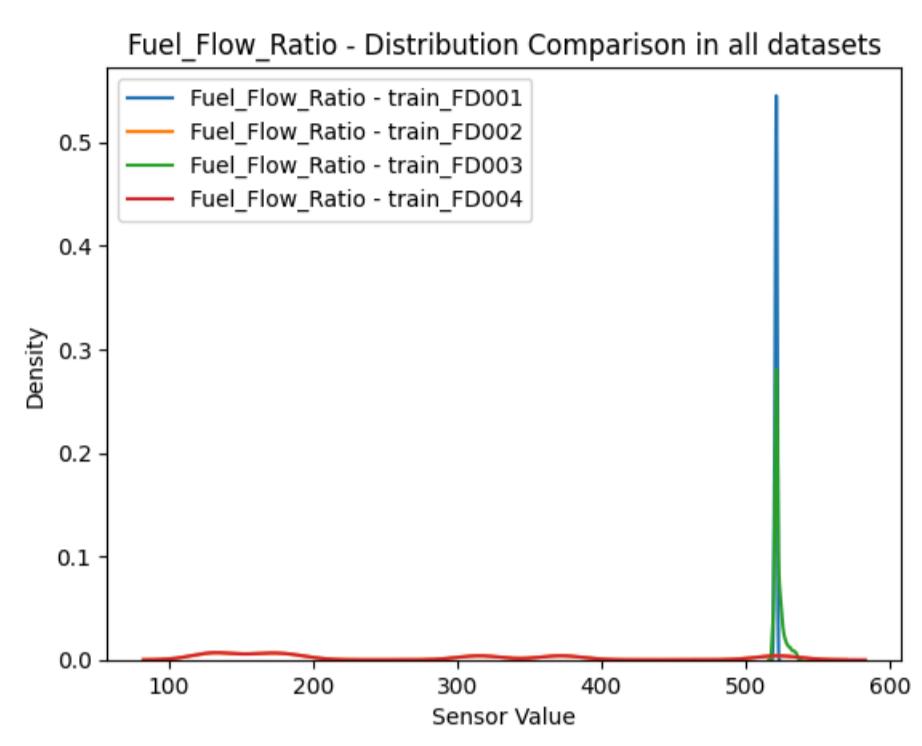
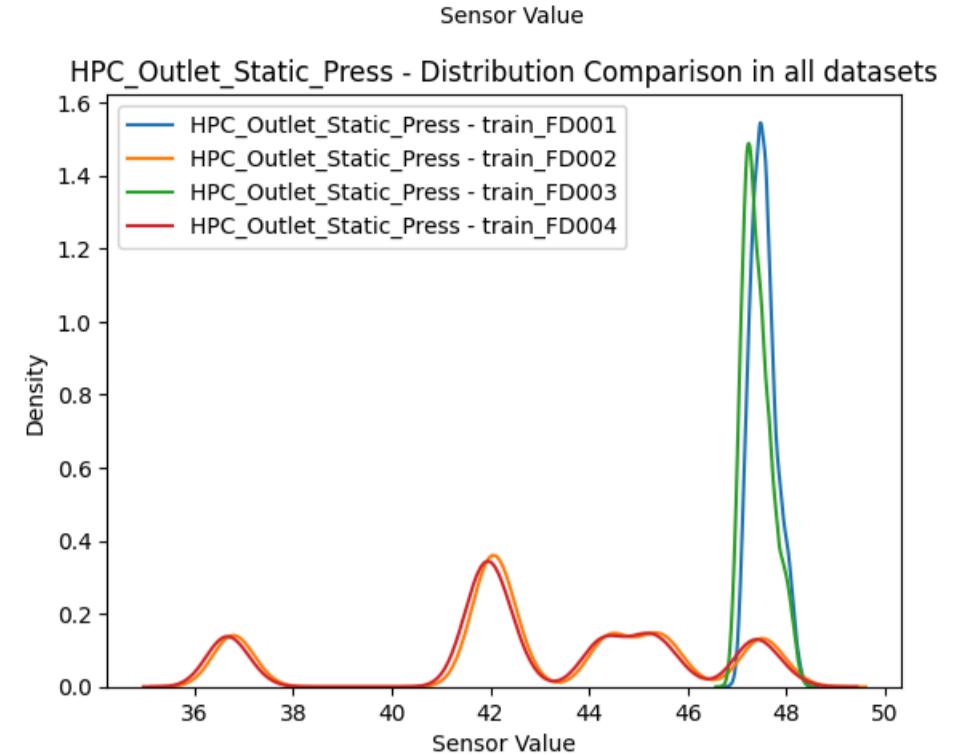
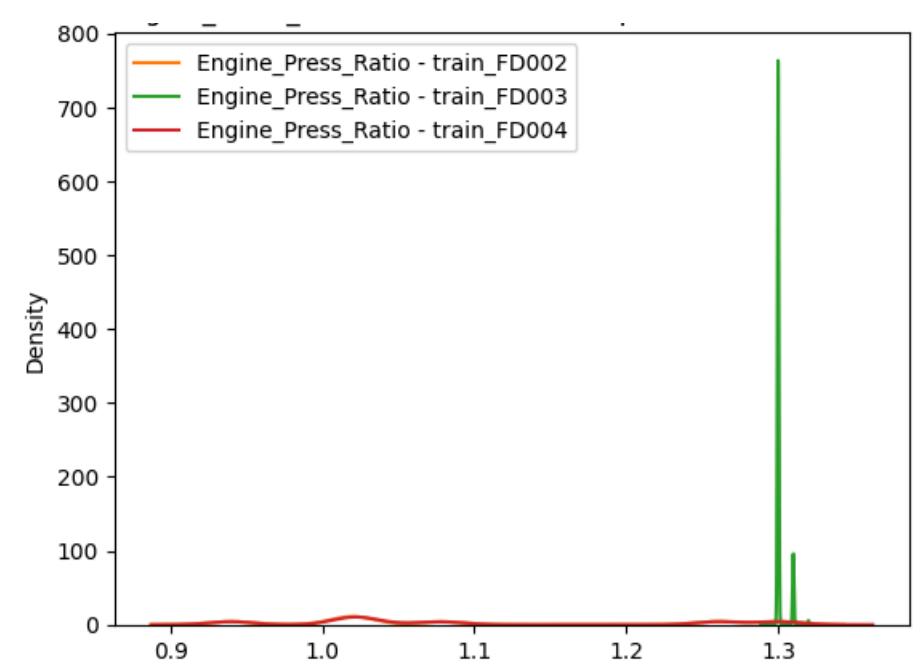
for sensor in sensors:
    for i, df in enumerate([train1, train2, train3, train4], start=1):
        sns.kdeplot(df[sensor], label=f"(sensor) - train_{i}0084")
    plt.title(f"(sensor) - Distribution Comparison in all datasets")
    plt.xlabel("Sensor Value")
    plt.ylabel("Density")
    plt.legend()
    plt.show()

# Summary comparison
for sensor in sensors:
    print(f"Summary for {sensor}: basic statistics:")
    for i, df in enumerate([train1, train2, train3, train4], start=1):
        print(f"({i}) mean={train_F0084[i]:.2f}, std={(df[sensor].std()):.2f}, min={(df[sensor].min()):.2f}, max={(df[sensor].max()):.2f}")
    print(f"({i+1}) mean={(df[sensor].mean()):.2f}, std={(df[sensor].std()):.2f}, min={(df[sensor].min()):.2f}, max={(df[sensor].max()):.2f}")
    print(f"({i+2}) mean={(df[sensor].mean()):.2f}, std={(df[sensor].std()):.2f}, min={(df[sensor].min()):.2f}, max={(df[sensor].max()):.2f}")
    print(f"({i+3}) mean={(df[sensor].mean()):.2f}, std={(df[sensor].std()):.2f}, min={(df[sensor].min()):.2f}, max={(df[sensor].max()):.2f}")

```









```

"Engine_Press_Ratio": {"min": 12.0, "max": 15.0}
}, "Packing_Station": {
    "Corrected_Core_Speed": {"min": -4800.0, "max": 4500.0},
    "Fuel_Air_Ratio": {"min": 0.02, "max": 0.03},
    "Req_Core_Speed": {"min": -3000.0, "max": 2500.0} # Approx
}
}

print("All mapping and ranges defined")
All mapping and ranges defined

Min - Max Normalization

# 57/85
def min_max_normalization(df, dataset_idx):
    df = df.copy()
    num_scaled = 0

    if dataset_idx in [0, 2, 4, 6]: # train1, train3, test1, test3 -> F0001/F0003
        UNIT_MACHINE_MAP = UNIT_MACHINE_MAP_01_03
    elif dataset_idx in [1, 5]: # train2, test2 -> F0002
        UNIT_MACHINE_MAP = UNIT_MACHINE_MAP_01_05
    elif dataset_idx in [3, 7]: # train4, test4 -> F0004
        UNIT_MACHINE_MAP = UNIT_MACHINE_MAP_04
    else:
        print(f"Invalid dataset_idx {dataset_idx}")
        return df

    for machine, unit_ids in UNIT_MACHINE_MAP.items():
        mask = df['unit'].isin(unit_ids)
        if mask.sum() == 0:
            print(f"Warning: No units for (machine) in dataset {dataset_idx}")
            continue
        dynamic_ranges = MACHINE_RANGES.get(machine, {})

        sensors = list(dynamic_ranges.keys()) + [col for col in df.columns if any(kw in col.lower() for kw in ['temp', 'press', 'speed', 'flow', 'ratio']) if col not in list(dynamic_ranges.keys()))[:2] # Ek 2 sensor max

        for sensor in sensors:
            if sensor in df.columns:
                sensor_data = df.loc[mask, sensor]

                q_min, q_max = sensor_data.quantile(0.001), sensor_data.quantile(0.999)
                if q_max - q_min < 0:
                    minv, maxv = q_min, q_max
                else:
                    minv, maxv = dynamic_ranges.get(sensor, {}).get('min', sensor_data.min()), dynamic_ranges.get(sensor, {}).get('max', sensor_data.max())

                scaled_series = (df.loc[mask, sensor] - minv) / (maxv - minv)

                scaled_series = np.clip(scaled_series, -1, 2)
                scaled_col_name = f'{sensor}_Scaled'
                df[scaled_col_name] = np.nan
                df.loc[mask, scaled_col_name] = scaled_series
                num_scaled += 1

    return df

datasets = [train1, train2, train3, train4, test1, test2, test3, test4]

for i, df in enumerate(datasets):
    datasets[i] = min_max_normalization(df, i)

# updating datasets
train1, train2, train3, train4, test1, test2, test3, test4 = datasets

print("Normalization successful")

# Train1 global plot
scaled_cols_train1 = [col for col in train1.columns if '_Scaled' in col]

plt.figure(figsize=(8,5))
plt.hist(train1[scaled_cols_train1[0]].dropna(), bins=50)
plt.title('Final Scaled Histogram (0-1, std>0.05 expected)')
plt.show()

Normalizing successful
Final Scaled Histogram (0-1, std>0.05 expected)


```

#### Removing Bias sensors and original columns

```

• original columns meaning before normalizations

# 59/85
print("Before Removing the original columns")
print(train1.shape)
train1.columns

Before Removing the original columns
(20631, 45)
['cycle', 'cycle1', 'gpcet', 'gpcet1', 'gpcet2', 'fan_inlet_Temp',
 'HPC_Outlet_Temp', 'HPC_Outlet_Temp', 'IPt_Outlet_Temp',
 'Fan_Inlet_Press', 'Bypass_Duct', 'HPC_Outlet_Press',
 'HPC_Outlet_Static_Press', 'HPC_Outlet_Static_Press', 'HPC_Outlet_Static_Press',
 'Corrected_Core_Speed', 'Bypass_Ratio', 'Fuel_Air_Ratio',
 'base_Rule', 'RUL', 'HPC_Cool_Air_F', 'HPC_Cool_Air_F',
 'IPt_Cool_Air_F', 'RUL', 'HPC_Outlet_Temp_Scaled',
 'HPC_Outlet_Press_Scaled', 'HPC_Cool_Air_F_Scaled',
 'IPt_Cool_Air_F_Scaled', 'Physical_Fan_Speed_Scaled',
 'Physical_Fan_Speed_Scaled', 'Physical_Core_Speed_Scaled',
 'HPC_Outlet_Static_Press_Scaled', 'HPC_Outlet_Static_Press_Scaled',
 'base_Rule', 'RUL', 'HPC_Outlet_Temp_Scaled',
 'IPt_Outlet_Temp_Scaled', 'Bypass_Duct_Scaled',
 'Fan_Inlet_Press_Scaled', 'Engine_Press_Ratio_Scaled',
 'Corrected_Core_Speed_Scaled', 'Fuel_Air_Ratio_Scaled',
 'Req_Core_Speed_Scaled'],
dtype='object')

# 60/85
# Bias sensor indecies (CMAPSS standard: column 0+1,5,6,10,16,18,19)
bias_cols = [5, 6, 10, 16, 18, 19] # keeping "cycle" 1 and "unit" 0 for RUL calculation

def clean_bias_and_original(df, dataset_idx):
    df = df.copy()
    df = df.drop(columns=[df.columns[i] for i in bias_cols])
    return df

scaled_cols = [col for col in df.columns if '_Scaled' in str(col)]
other_cols = ['unit', 'cycle']
df = df[other_cols + scaled_cols]

print(f"Cleaned dataset {dataset_idx}: shape: {df.shape}, scaled cols: {len(scaled_cols)}")
return df

datasets = [train1, train2, train3, train4, test1, test2, test3, test4]
for i, df in enumerate(datasets):
    datasets[i] = clean_bias_and_original(df, i)
train1, train2, train3, train4, test1, test2, test3, test4 = datasets

Cleaned dataset 0: Shape: (20631, 20), Scaled cols: 18
Cleaned dataset 1: Shape: (20631, 20), Scaled cols: 18
Cleaned dataset 2: Shape: (20728, 20), Scaled cols: 18
Cleaned dataset 3: Shape: (61249, 20), Scaled cols: 18
Cleaned dataset 4: Shape: (13096, 20), Scaled cols: 18
Cleaned dataset 5: Shape: (13096, 20), Scaled cols: 18
Cleaned dataset 6: Shape: (165596, 20), Scaled cols: 18
Cleaned dataset 7: Shape: (41214, 20), Scaled cols: 18

# 61/85
print("After Removing the original columns")
print(train1.shape)
train1.columns

After Removing the original columns
(20631, 20)
['cycle', 'HPC_Outlet_Temp_Scaled', 'HPC_Outlet_Press_Scaled',
 'IPt_Outlet_Temp_Scaled', 'Fan_Inlet_Press_Scaled',
 'Physical_Core_Speed_Scaled', 'HPC_Outlet_Static_Press_Scaled',
 'Corrected_Core_Speed_Scaled', 'Bypass_Duct_Scaled',
 'IPt_Outlet_Temp_Scaled', 'HPC_Outlet_Static_Press_Scaled',
 'Fan_Inlet_Press_Scaled', 'Engine_Press_Ratio_Scaled',
 'Corrected_Core_Speed_Scaled', 'Fuel_Air_Ratio_Scaled',
 'Req_Core_Speed_Scaled'],
dtype='object')

```

#### RUL Calculation & Adding

```

# 63/85
def compute_rul(df, max_rul=125):
    df = df.copy()
    df['RUL'] = np.nan
    for unit in df['unit'].unique():
        max_cycle = df[df['unit'] == unit]['cycle'].max()
        df[(df['unit'] == unit) & (df['RUL'] == max_cycle)] = max_cycle - df['cycle']
        df[(df['unit'] == unit) & (df['RUL'] == max_cycle)] = np.clip(df['RUL'], 0, max_rul)
    return df

train1 = compute_rul(train1)
train2 = compute_rul(train2)
train3 = compute_rul(train3)
train4 = compute_rul(train4)
print("Train datasets RUL calculation completed and added\n")

def compute_rul_test(df, rul_df):
    df = df.copy()
    rul = df['RUL'].values[1:] # rul for 1, rul in enumerate(rul_df['RUL'])) # unit 1-based
    df['RUL'] = df['unit'].map(rul_df['RUL'])
    return df

test1 = compute_rul_test(test1, rul1)
test2 = compute_rul_test(test2, rul2)
test3 = compute_rul_test(test3, rul3)
test4 = compute_rul_test(test4, rul4)
print("Test datasets RUL calculation completed and added")
Train datasets RUL calculation completed and added
Test datasets RUL calculation completed and added

# 64/85
print("Train RUL Stats (Expected: mean=100, min=0, max=125):")
print(train1['RUL'].describe())
print("Train RUL per unit sample (unique=1, variable):")
print(train1.groupby('unit')['RUL'].agg(['min', 'max', 'nunique']).head(5))

print("Test1 RUL Stats (Expected: mean=128, constant per unit):")
print(test1['RUL'].describe())
print("Test1 RUL per unit sample (unique=1, constant):")
print(test1.groupby('unit')['RUL'].agg(['first', 'nunique']).head(5)) # first=constant RUL
print("Test1 RUL unique count (Expected: 100 unit):", test1['RUL'].nunique())

print("\nAny NaN in RUL Train1:", train1['RUL'].isna().sum(), "Test1:", test1['RUL'].isna().sum())
Train RUL Stats (Expected: mean=100, min=0, max=125):
count    20631.000000
mean     86.829285
std      10.873390
min      0.000000
25%     101.000000
50%     101.000000
75%     125.000000
max    145.000000
Name: RUL, dtype: float64

Train RUL per unit sample (unique=1, variable):
unit
1    0.0   125.0   126
2    0.0   125.0   126
3    0.0   125.0   126
4    0.0   125.0   126
5    0.0   125.0   126

Test1 RUL Stats (Expected: mean=128, constant per unit):
count    13096.000000
mean     65.446195
std      42.354423
min     1.000000
25%     21.000000
50%     50.000000
75%     97.000000
max    145.000000
Name: RUL, dtype: float64

Test1 RUL per unit sample (unique=1, constant):
unit
1    0.0   128.0   126
2    0.0   128.0   126
3    0.0   128.0   126
4    0.0   128.0   126
5    0.0   128.0   126

Test1 RUL unique count (Expected: 100 unit): 71
Any NaN in RUL Train1: 0 Test1: 0

```

#### Window Based Feature Aggregation

```
# 66/85

print("Before window - RUL sample:", train1['RUL'].head())
print("Before - Shape:", train1.shape)
print("Before - Index sample:", train1.index[:5].tolist())
Before window - RUL sample: 0  125.0
1  125.0
2  125.0
3  125.0
Name: RUL, dtype: float64
Before - Shape: (20631, 21)
Before - Index sample: [0, 1, 2, 3, 4]

# 67/85
# (Updated 22/10/25 for Overfit Fix: Shifted Window)

def window_features(df, window_size=30, shift=1):
    df = df.copy().sort_values(['unit', 'cycle']).reset_index(drop=True)
    rul_col = None
    if 'RUL' in df.columns:
        rul_col = df['RUL'].copy()
    feature_cols = [col for col in df.columns if '_Scaled' in str(col)]

    for col in feature_cols:
        rolling_mean = df.groupby('unit')[col].rolling(window=window_size, min_periods=1).mean().shift(shift).reset_index(0, drop=True).values
        rolling_std = df.groupby('unit')[col].rolling(window=window_size, min_periods=1).std().shift(shift).reset_index(0, drop=True).values
        rolling_var = df.groupby('unit')[col].rolling(window=window_size, min_periods=1).var().shift(shift).reset_index(0, drop=True).values
        df[f'{col}_mean'] = rolling_mean
        df[f'{col}_std'] = rolling_std
        df[f'{col}_var'] = rolling_var

    # Groupby fillna for stat
    agg_cols = [(f'{col}_stat') for col in feature_cols for stat in ['mean', 'std', 'var']]
    df[agg_cols] = df.groupby('unit')[agg_cols].ffill()

    # Fillna(0) for any remaining non-RUL Nans
    non_rul_col = [col for col in df.columns if col != 'RUL']
    df[non_rul_col] = df[non_rul_col].fillna(0)

    if rul_col is not None:
        df['RUL'] = rul_col

    if 'RUL' in df.columns:
        print(f'{df["RUL"].isna().sum()} Nans in RUL after window')
    return df

train1 = window_features(train1)
train2 = window_features(train2)
train3 = window_features(train3)
train4 = window_features(train4)

print("Window features applied to all train datasets")

test1 = window_features(test1)
test2 = window_features(test2)
test3 = window_features(test3)
test4 = window_features(test4)

print("Window features applied to all test datasets")

0 Nans in RUL after window
Window features applied to all train datasets
0 Nans in RUL after window
Window features applied to all test datasets

# 68/85

print("After - RUL sample:", train1['RUL'].head())
print("After - Shape:", train1.shape)
print("After - Index sample:", train1.index[:5].tolist())
print("RUL match?", np.allclose(train1['RUL'].values, train1['RUL'].values))
print("Any NAn in RUL after?", train1['RUL'].isna().sum())

After - RUL sample: 0  125.0
1  125.0
2  125.0
3  125.0
4  125.0
Name: RUL, dtype: float64
After - Shape: (20631, 75)
After - Index sample: [0, 1, 2, 3, 4]
RUL match? True
Any NAn in RUL after? 0
```

#### Merging All Train Datasets

```
# 70/85
# (Updated 22/10/25 for Overfit Fix: removed .fillna)
full_merged = pd.concat([train1, train2, train3, train4], ignore_index=True)
print("Full Merged shape:", full_merged.shape)
print("Full RUL stats:", full_merged['RUL'].describe())
print("Any NAn in features?", full_merged.drop(['RUL', 'cycle', 'unit'], axis=1).isna().sum())

Full Merged shape: (160559, 75)
full_rul = full_merged['RUL']
count      160559.000000
mean       98.182629
std        16.151556
min       0.000000
25%      16.000000
50%      118.000000
75%      125.000000
max      125.000000
name: RUL, dtype: float64
Any NAn in features? 0
```

#### Synthetic Data Integration & Data Balance

##### Synthetic Data Generation

```
# 73/85
# (Updated 22/10/25 for Overfit Fix: Moved from after Feature selection (72/85) to before 73/85)
# (Updated 22/10/25 for Overfit Fix: Col Align + Post-Merge /fill for Nuhw)
# (Updated 23/10/25 for Overfit Fix: raw only baseline filter // no reindex // single agg layer)

n_synthetic = 50000
new_units = n_synthetic * 100
raw_units = np.random.randint(100, 150, 50) # 50 unique units

# Raw: RUL is random, scaled but exclude agg suffixes (no opset, bias dropped)
raw_SUFFIXES = ['_mean', '_std', '_var']
raw_cols = [col for col in full_merged.columns if '_Scaled' in col and not any(suf in col for suf in raw_SUFFIXES)]
print(f"Raw scaled cols for baseline: {len(raw_cols)} ({raw_cols[:3]}...)" ) # ~18

baseline_raw = full_merged[raw_cols].mean()
drift_rate = 0.1 * full_merged[raw_cols].std()
noise_rnd = np.random.normal(0, 1, len(raw_cols))

for i in range(new_units):
    raw_col = raw_cols[i % len(raw_cols)]
    raw_val = baseline_raw[raw_col] + drift_rate[raw_col] * noise_rnd[i]
    raw_val += noise_rnd[i] * noise_rnd[i]
```