```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Mounted at /content/drive

```
1 import sys
2 sys.path.append( '/content/drive/MyDrive/CS 480 Kaggle Competition/PlantTraits2024/src')
```

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import os
```

```
1 data_path = os.getcwd()
2 df_train = pd.read_csv('/content/drive/MyDrive/CS 480 Kaggle Competition/data/train.csv')
3 df_train['path'] = '/content/drive/MyDrive/CS 480 Kaggle Competition/data/train.csv' + df_train['id'].astype(str) + '.jpeg'
4 df_test = pd.read_csv('/content/drive/MyDrive/CS 480 Kaggle Competition/data/test.csv')
5 df_test['path'] = '/content/drive/MyDrive/CS 480 Kaggle Competition/data/test.csv' + df_test['id'].astype(str) + '.jpeg'
```
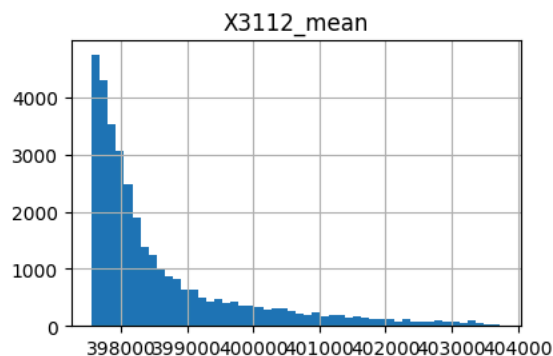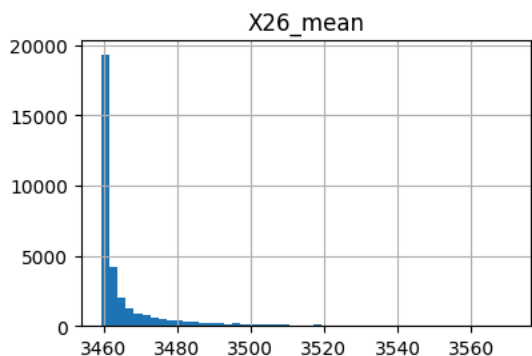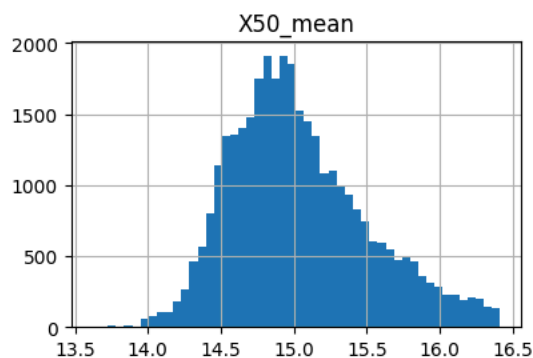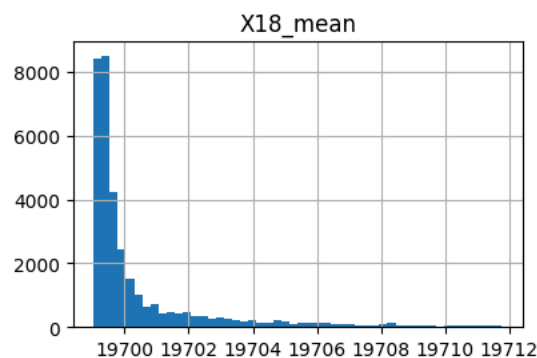
```
1 df_train
```

| nality | WORLDCLIM_BIO4_temperature_seasonality | WORLDCLIM_BIO7_temperature_annual_range |
|---|---|---|
| 210766 | 161.457764 | 13.886666 |
| 906487 | 178.745422 | 19.846668 |
| 545128 | 292.781219 | 23.486668 |
| 563957 | 211.065521 | 16.768000 |
| 409706 | 36.499138 | 10.257143 |
| ... | ... | ... |
| 970898 | 75.369301 | 12.087244 |
| 597702 | 120.009247 | 14.226222 |
| 789906 | 473.979675 | 26.604889 |
| 718536 | 182.917358 | 22.998470 |
| 389532 | 1090.754761 | 41.099998 |

```
1 #all columns must be identical to be consider the same species
2 trait_columns = ['X4_mean', 'X11_mean', 'X18_mean', 'X50_mean', 'X26_mean', 'X3112_mean']
3 aux_columns = list(
4           map(lambda x: x.replace("mean", "sd"), trait_columns)
5         )
6 metadata_cols = df_train.drop(
7           columns=["id", "path"] + trait_columns
8         ).columns
```

```
1 for col in trait_columns:
2    upper_quantile = df_train[col].quantile(0.98)
3    df_train = df_train[(df_train[col] < upper_quantile)]
4    df_train = df_train[(df_train[col] > 0)]
```
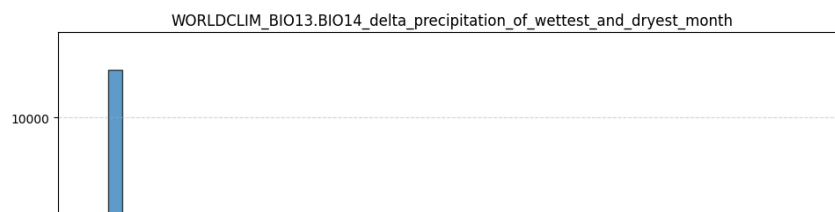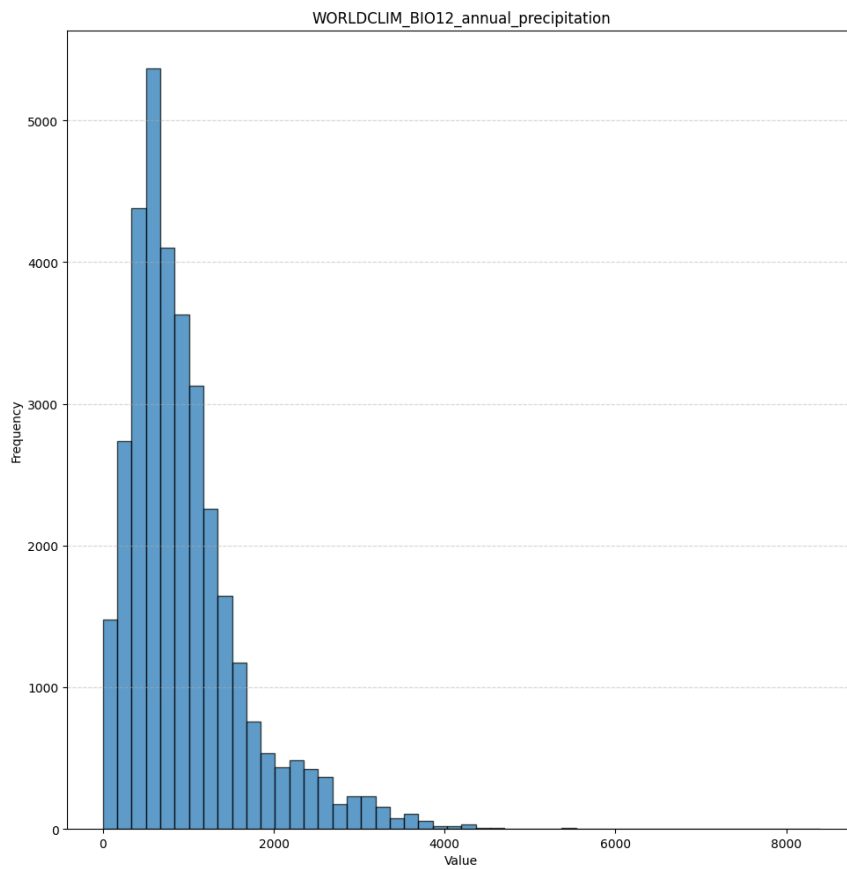
## Data Visualization of variables to be predicted
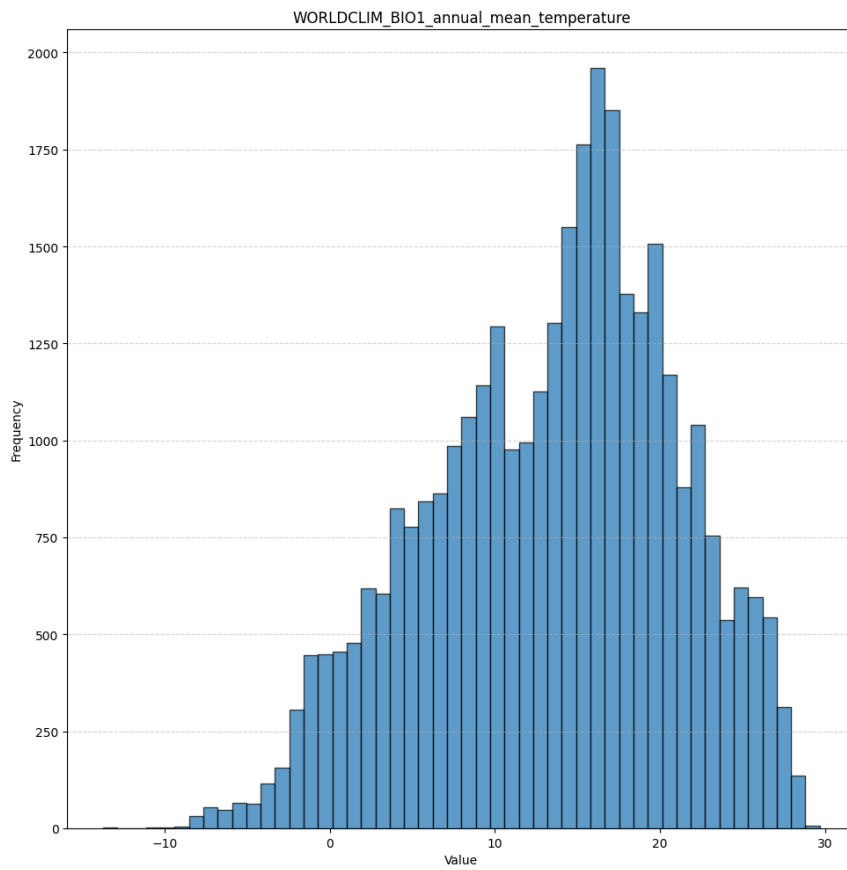
```
1 df_train[trait_columns].hist(bins=50, figsize=(10, 10))
2 plt.show()
3
```
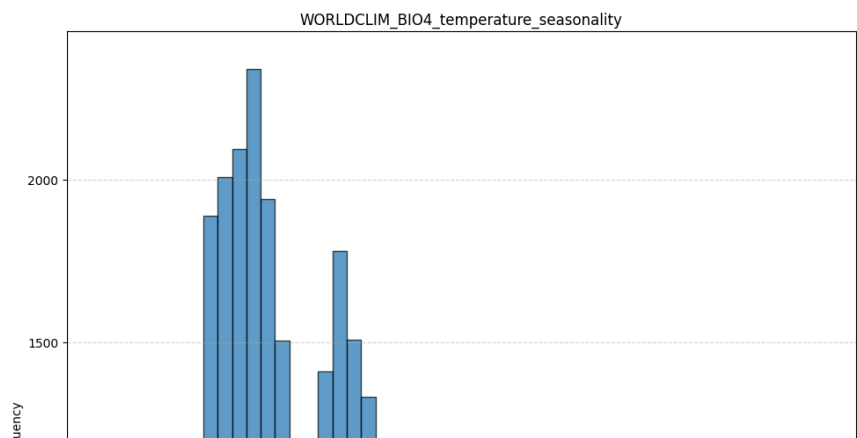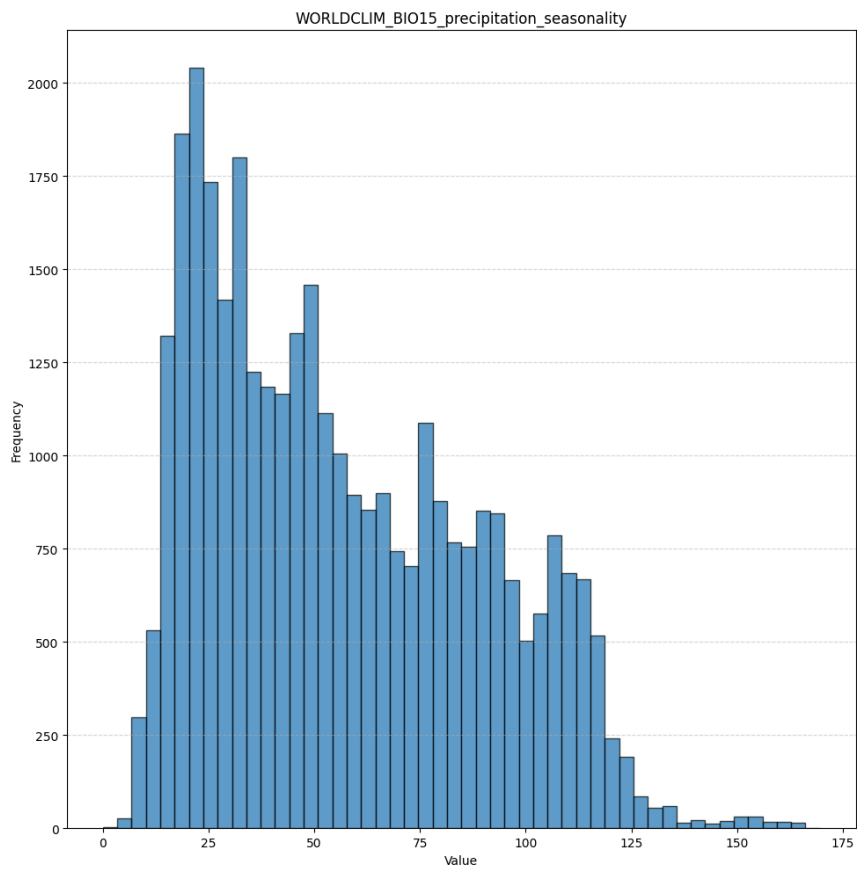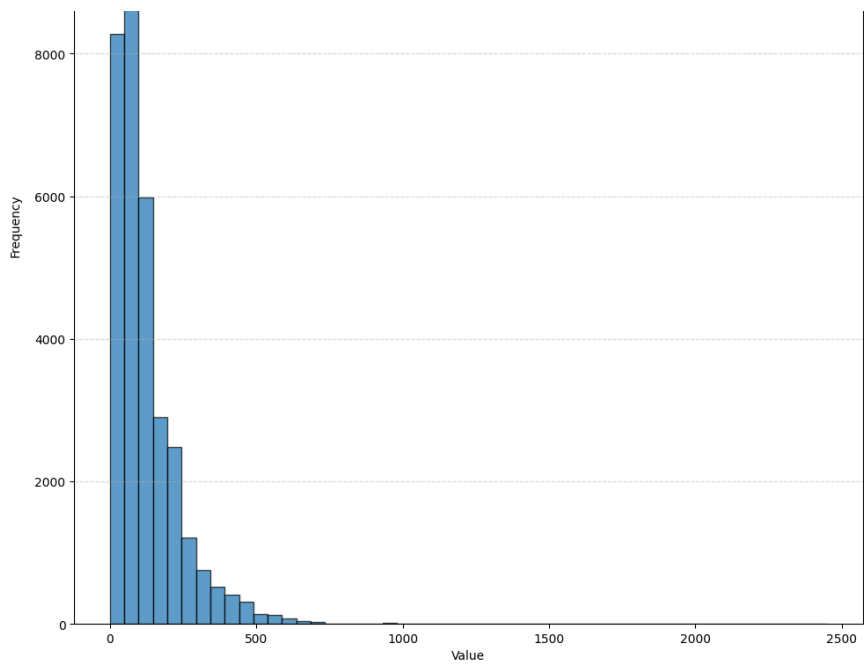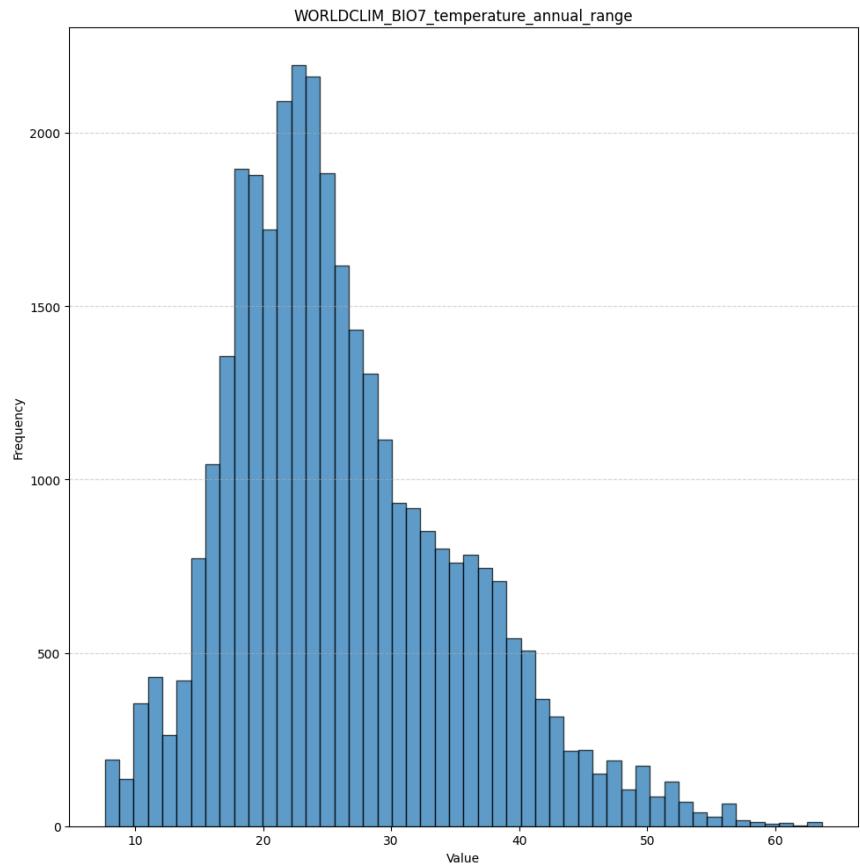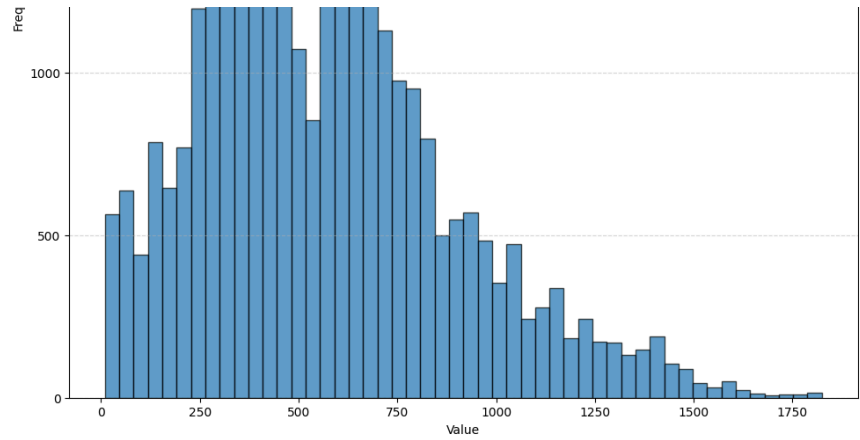


```
1 def visualize(df_train, colQuery):
2   climate_cols = df_train.filter(like=colQuery)
3   fig, axs = plt.subplots(nrows=len(climate_cols.columns), ncols=1, figsize=(10, 70))
4
5   # loop through each column and create a histogram
6   for i, col in enumerate(climate_cols.columns):
7       axs[i].hist(climate_cols[col], bins=50, alpha=0.7, edgecolor = 'black')
8       axs[i].set_title(col)
9       axs[i].set_xlabel('Value')
10      axs[i].set_ylabel('Frequency')
11      axs[i].grid(True, axis='y', linestyle='--', alpha=0.5)
12
13  # adjust the layout to avoid overlapping titles
14  plt.tight_layout(rect=[0, 0.03, 1, 0.95])
15
16
17  # show the plot
18  plt.show()
```

## ⌄ WORLDCLIM_BIO* Data Visualization

```
1 visualize(df_train, 'WORLDCLIM_BIO')
```

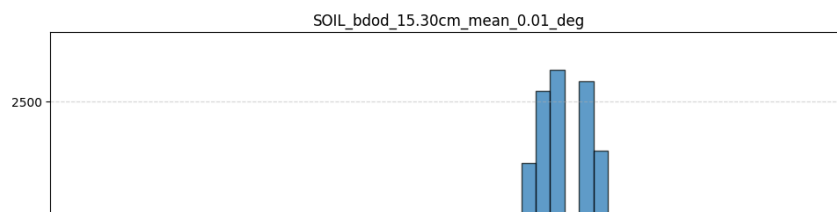### WORLDCLIM_BIO1_annual_mean_temperature



### WORLDCLIM_BIO12_annual_precipitation



### WORLDCLIM_BIO13.BIO14_delta_precipitation_of_wettest_and_dryest_month

WORLDCLIM_BIO15_precipitation_seasonality



WORLDCLIM_BIO4_temperature_seasonality

WORLDCLIM_BIO7_temperature_annual_range

## ⌄ SOIL_bdod* Data Visualization

```
1 visualize(df_train, 'SOIL_bdod')
```

SOIL_bdod_0.5cm_mean_0.01_deg



SOIL_bdod_100.200cm_mean_0.01_deg



SOIL_bdod_15.30cm_mean_0.01_deg

SOIL_bdod_30.60cm_mean_0.01_deg



SOIL_bdod_5.15cm_mean_0.01_deg

SOIL_bdod_60.100cm_mean_0.01_deg

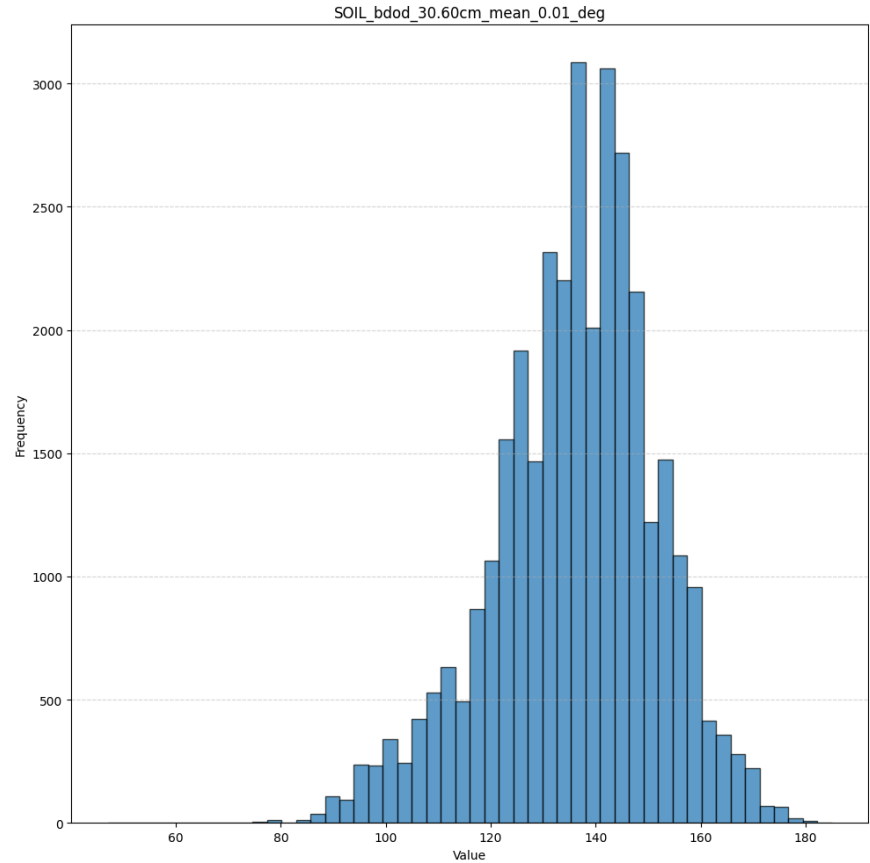## ⌄ VOD Data Visualization

```
1 visualize(df_train, 'VOD')
```

VOD_C_2002_2018_multiyear_mean_m01



VOD_C_2002_2018_multiyear_mean_m02



VOD_C_2002_2018_multiyear_mean_m03



VOD_C_2002_2018_multiyear_mean_m04



VOD_C_2002_2018_multiyear_mean_m05



VOD_C_2002_2018_multiyear_mean_m06



VOD_C_2002_2018_multiyear_mean_m07



VOD_C_2002_2018_multiyear_mean_m08



VOD_C_2002_2018_multiyear_mean_m09



VOD_C_2002_2018_multiyear_mean_m10



VOD_C_2002_2018_multiyear_mean_m11



VOD_C_2002_2018_multiyear_mean_m12



VOD_Ku_1987_2017_multiyear_mean_m01



VOD_Ku_1987_2017_multiyear_mean_m02