# Exploratory Data Analysis (EDA) on RSNA 2022 HDF5 Subset (Colab Notebook)

```
1 # Install Kaggle API
2 !pip install -q kaggle
```

```
 1 # Import libraries
 2 import os
 3 from google.colab import files
 4 import sys
 5 import h5py
 6 import numpy as np
 7 import matplotlib.pyplot as plt
 8 import matplotlib.patches as patches
 9 import random
10 import seaborn as sns
11 import pandas as pd
12 from tqdm.auto import tqdm
```

```
 1 # Check if the API key already exists
 2 if not os.path.exists('/root/.kaggle/kaggle.json'):
 3     print('Please upload your personal kaggle.json file:')
 4
 5     # Prompt for upload
 6     uploaded = files.upload()
 7
 8     # Move the file to the correct directory
 9     !mkdir -p ~/.kaggle
10     !cp kaggle.json ~/.kaggle/
11     !chmod 600 ~/.kaggle/kaggle.json
12     print('\nKaggle API key configured successfully.')
13 else:
14     print('Kaggle API key is already configured.')
```

```
Please upload your personal kaggle.json file:
Choose Files  No file chosen          Upload widget is only available when the cell has been executed in the current browser session.
Please rerun this cell to enable.
Saving kaggle.json to kaggle.json

Kaggle API key configured successfully.
```

```
1 # Replace with your dataset id here
2 kaggle_dataset_id = 'andymalinsky/rsna-2022-hdf5-subset'
3
4 print('Downloading dataset...')
5 !kaggle datasets download -d {kaggle_dataset_id}
```

```
Downloading dataset...
Dataset URL: https://www.kaggle.com/datasets/andymalinsky/rsna-2022-hdf5-subset
License(s): CC0-1.0
Downloading rsna-2022-hdf5-subset.zip to /content
 99% 3.49G/3.51G [00:53<00:00, 275MB/s]
100% 3.51G/3.51G [00:53<00:00, 70.7MB/s]
```

```
1 # Unzip and move to the local Colab storage
2 zip_name = kaggle_dataset_id.split('/')[-1] + '.zip'
3 !unzip -q '{zip_name}'
4
5 # This is the path to the data
6 data_path = 'fracture_dataset_subset.h5'
```

```
1 # Open and view the file
2 with h5py.File(data_path, 'r') as f:
3     # View key list, essentially our 'columns'
4     print(f'Keys: {list(f.keys())}')
5
6     # View column shapes
7     print('\n--- Details ---')
8     for key in f.keys():
9         dataset = f[key]
10        print(f'\nKey: {key}')
11        print(f'Shape: {dataset.shape}')
12        print(f'Dtype: {dataset.dtype}')
```

```
Keys: ['SliceNumber', 'StudyInstanceUID', 'bboxes', 'images', 'labels', 'split']

--- Details ---

Key: SliceNumber
Shape: (28812,)
Dtype: object

Key: StudyInstanceUID
Shape: (28812,)
Dtype: object

Key: bboxes
Shape: (28812, 10, 4)
Dtype: float32

Key: images
Shape: (28812, 256, 256)
Dtype: float32

Key: labels
Shape: (28812,)
Dtype: int8

Key: split
Shape: (28812,)
Dtype: int8
```

```
1 # Define visualizer function
2 def visualize_samples(num_samples, label):
3     with h5py.File(data_path, 'r') as hf:
4         # Load all labels and find the indices for the desired sample type
5         all_labels = hf['labels'][:]
6         indices_to_sample_from = np.where(all_labels == label)[0]
7         sample_indices = random.sample(list(indices_to_sample_from), num_samples)
8         fig, axes = plt.subplots(1, num_samples, figsize=(10, 10))
9
10        for ax, idx in zip(axes, sample_indices):
11            # Load the data for this one sample
12            image = hf['images'][idx]
13            bboxes = hf['bboxes'][idx]
14            uid = hf['StudyInstanceUID'][idx]
15            slice_num = hf['SliceNumber'][idx]
16            label = hf['labels'][idx]
17
18            # Filter out the "-1" padding on the bounding boxes
19            bboxes = [box for box in bboxes if box[0] != -1.0]
20
21            # Plot the image
22            ax.imshow(image, cmap='bone')
23            # Set title
24            ax.set_title(f'Negative Sample (Label: {label}) \nSlice: {slice_num} | (Index: {idx})')
25
26            # Plot each bounding box
27            for box in bboxes:
28                x, y, w, h = box # the box value contains its coordinates and dimensions
29
30                # Draw the bounding box on the image
31                rect = patches.Rectangle(
```

```
32                    (x, y), w, h,
33                    linewidth=2,
34                    edgecolor='r',
35                    facecolor='none'
36                )
37                ax.add_patch(rect)
38
39        plt.tight_layout()
40        plt.show()
```
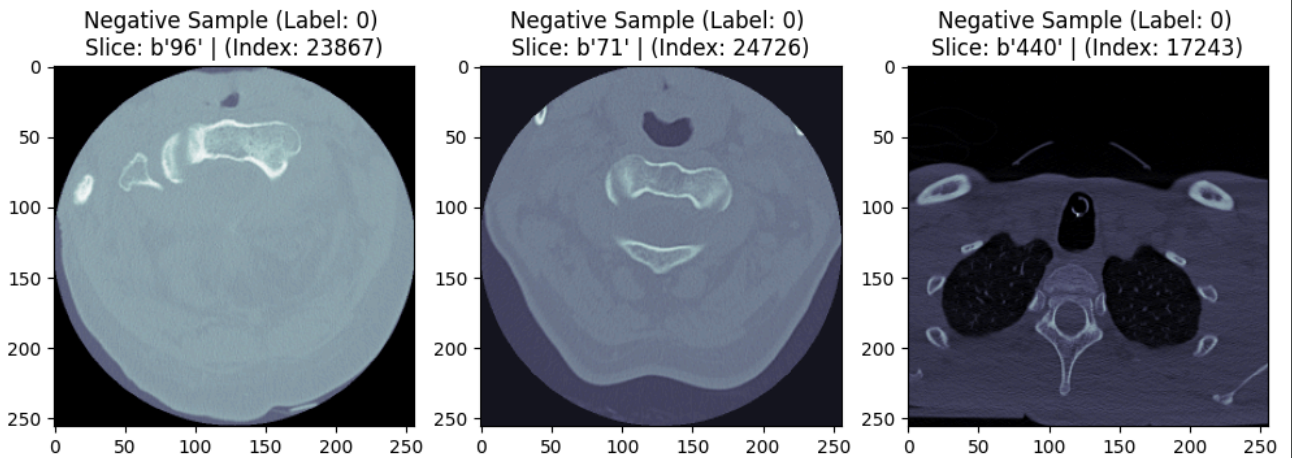
```
1 # View 3 sample negative images
2 visualize_samples(3, 0)
```
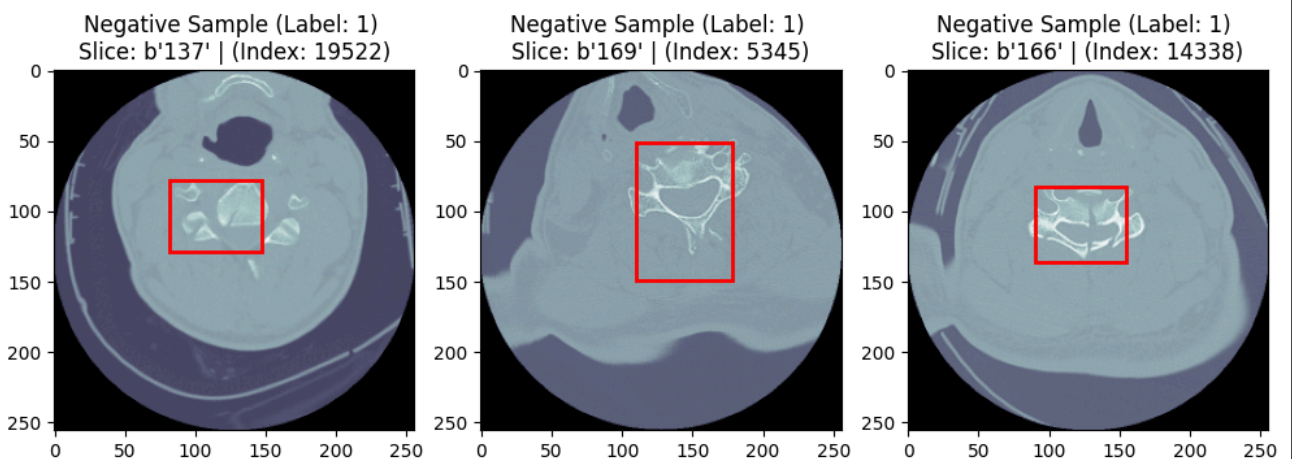


```
1 # View 3 sample positive images
2 visualize_samples(3, 1)
```
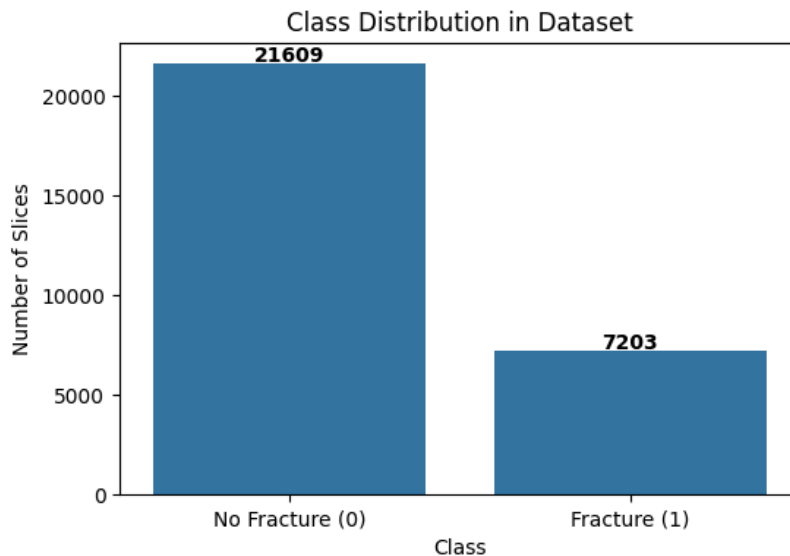


```
1 # Check class balance
2
3 # Get labels
4 with h5py.File(data_path, 'r') as f:
5     labels = f['labels'][:]
6
7 # Calculate counts
8 unique, counts = np.unique(labels, return_counts=True)
9 class_names = ['No Fracture (0)', 'Fracture (1)']
10
11 # Plot
12 plt.figure(figsize=(6, 4))
```

```
13 sns.barplot(x=class_names, y=counts)
14 plt.title("Class Distribution in Dataset")
15 plt.ylabel("Number of Slices")
16 plt.xlabel("Class")
17
18 # Add count labels on top
19 for i, v in enumerate(counts):
20     plt.text(i, v + 150, str(v), ha='center', fontweight='bold')
21
22 plt.show()
23 print(f'Ratio: 1 Positive (Fracture) for every {int(counts[0]/counts[1])} Negative (No Fracture)')
```



Class Distribution in Dataset

```
Ratio: 1 Positive (Fracture) for every 3 Negative (No Fracture)
```

```
 1 # Get labels with splits
 2 with h5py.File(data_path, 'r') as f:
 3     splits = f['split'][:] # 0=Train, 1=Val, 2=Test
 4     labels = f['labels'][:] # 0=Negative, 1=Positive
 5
 6 # Create dataFrame
 7 eda_df = pd.DataFrame({
 8     'split_id': splits,
 9     'label_id': labels
10 })
11
12 # Map id to split/label names
13 split_map = {0: 'Train', 1: 'Validation', 2: 'Test'}
14 label_map = {0: 'No Fracture', 1: 'Fracture'}
15 eda_df['split'] = eda_df['split_id'].map(split_map)
16 eda_df['label'] = eda_df['label_id'].map(label_map)
17
18 # Counts by split and label
19 cnt = eda_df.groupby(['split', 'label']).size().reset_index(name='count')
20
21 # Reorder for consistent display (Train -> Val -> Test)
22 cnt['split'] = pd.Categorical(cnt['split'], categories=['Train', 'Validation', 'Test'], ordered=True)
23 cnt = cnt.sort_values('split')
24 print('\x1B[4m' + 'Sample Counts by Split and Class' + '\x1B[0m')
25 display(cnt)
26
27 # Check if ratio is consistent across splits
28 print('\n\x1B[4m' + 'Fracture Rate per Split' + '\x1B[0m')
29 for split_name in ['Train', 'Validation', 'Test']:
30     subset = eda_df[eda_df['split'] == split_name]
31     pos_rate = subset['label_id'].mean()
32     print(f'{split_name}: {pos_rate:.1%} Positive')
```

Sample Counts by Split and Class

|   | split | label | count |
|---|-------|-------|-------|
| 2 | Train | Fracture | 5186 |
| 3 | Train | No Fracture | 15558 |
| 4 | Validation | Fracture | 1045 |
| 5 | Validation | No Fracture | 3135 |
| 0 | Test | Fracture | 972 |
| 1 | Test | No Fracture | 2916 |

Fracture Rate per Split
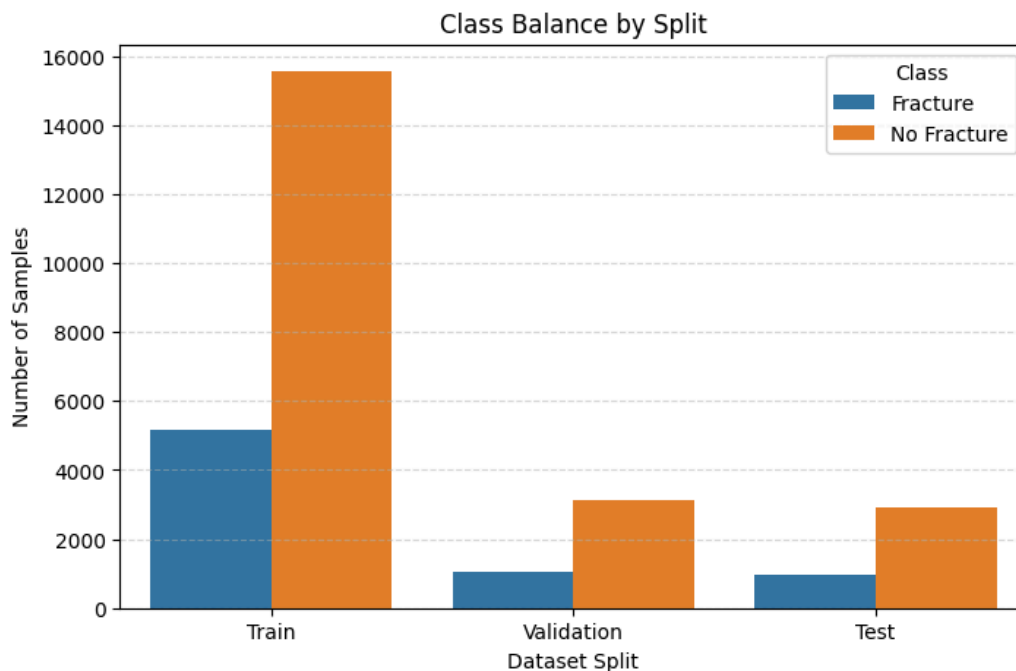Train: 25.0% Positive
Validation: 25.0% Positive
Test: 25.0% Positive

```python
1 # Plot class balance by split
2 plt.figure(figsize=(8, 5))
3 sns.barplot(data=cnt, x='split', y='count', hue='label')
4 plt.title('Class Balance by Split')
5 plt.xlabel('Dataset Split')
6 plt.ylabel('Number of Samples')
7 plt.legend(title='Class')
8 plt.grid(axis='y', linestyle='--', alpha=0.5)
9 plt.show()
```
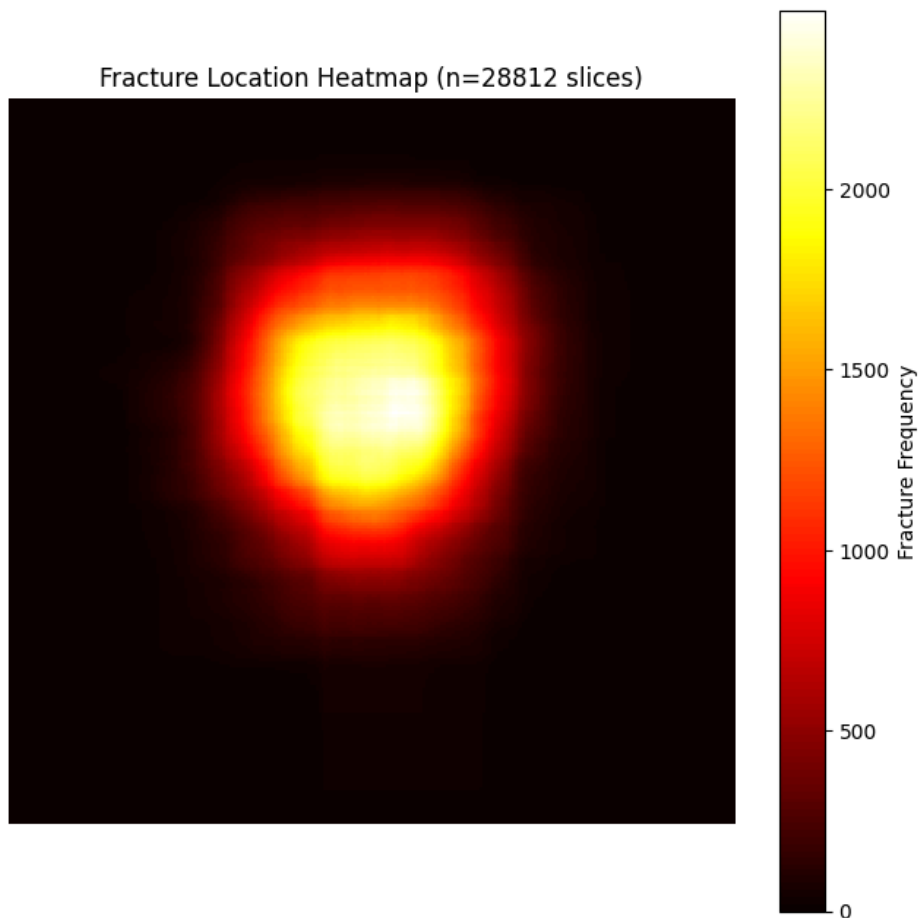


```python
1 # Plot heatmap of fracture locations
2 # This shows us where fractures are typically found in the images
3
4 # Initialize a blank heatmap
5 heatmap = np.zeros((256, 256))
6
7 # Get all bounding boxes
8 with h5py.File(data_path, 'r') as f:
9     all_bboxes = f['bboxes'][:]
10
11     # Loop through all samples
12     for i in range(len(all_bboxes)):
13         # Get bounding box for this sample
14         bbox = all_bboxes[i]
15
```

```
16        # Filter out -1 padding
17        valid_box = bbox[bbox[:, 0] != -1]
18
19        for box in valid_box:
20            # Box format in H5 is [x, y, w, h] (pixels)
21            # We need to make sure it's within 256x256 bounds
22            x, y, w, h = box.astype(int)
23
24            # Simple check to ensure indices are valid
25            x = max(0, x)
26            y = max(0, y)
27            w = min(256 - x, w)
28            h = min(256 - y, h)
29
30            # Add 1 to the heatmap region corresponding to the box
31            heatmap[y:y+h, x:x+w] += 1
32
33 # Plot the heatmap
34 plt.figure(figsize=(8, 8))
35 plt.imshow(heatmap, cmap='hot', interpolation='nearest')
36 plt.colorbar(label='Fracture Frequency')
37 plt.title(f'Fracture Location Heatmap (n={len(all_bboxes)} slices)')
38 plt.axis('off')
39 plt.show()
```



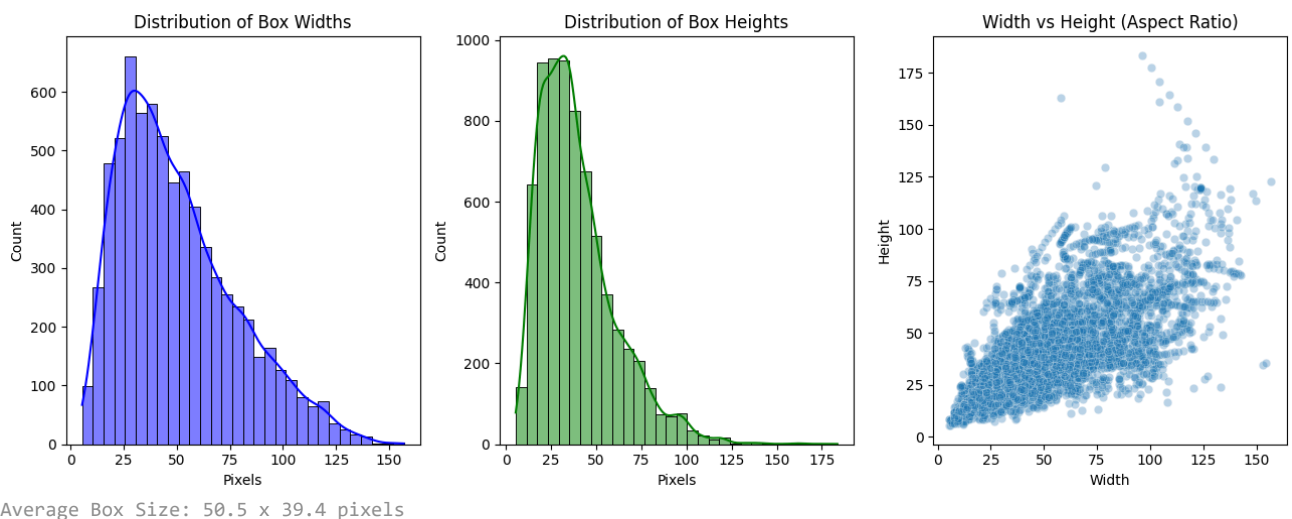Fracture Location Heatmap (n=28812 slices)

```
1 # Explore bounding box sizes
2 box_widths = []
3 box_heights = []
4 box_areas = []
5
6 # Open file and get all bounding boxes
7 with h5py.File(data_path, 'r') as f:
8     all_bboxes = f['bboxes'][:]
9
10     for boxes in all_bboxes:
```

```
11        # Filter out -1 padding
12        valid_boxes = boxes[boxes[:, 0] != -1]
13
14        for box in valid_boxes:
15            w, h = box[2], box[3]
16            box_widths.append(w)
17            box_heights.append(h)
18            box_areas.append(w * h)
19
20 # Plot histograms
21 fig, ax = plt.subplots(1, 3, figsize=(13, 5))
22
23 sns.histplot(box_widths, bins=30, ax=ax[0], color='blue', kde=True)
24 ax[0].set_title('Distribution of Box Widths')
25 ax[0].set_xlabel('Pixels')
26
27 sns.histplot(box_heights, bins=30, ax=ax[1], color='green', kde=True)
28 ax[1].set_title('Distribution of Box Heights')
29 ax[1].set_xlabel('Pixels')
30
31 sns.scatterplot(x=box_widths, y=box_heights, ax=ax[2], alpha=0.3)
32 ax[2].set_title('Width vs Height (Aspect Ratio)')
33 ax[2].set_xlabel('Width')
34 ax[2].set_ylabel('Height')
35
36 plt.tight_layout()
37 plt.show()
38
39 print(f'Average Box Size: {np.mean(box_widths):.1f} x {np.mean(box_heights):.1f} pixels')
```



```
Average Box Size: 50.5 x 39.4 pixels
```

```
1 # Check pixel intensity distribution
2
3 # Sample 1000 random images just to check
4 num_samples_to_check = 1000
5
6 with h5py.File(data_path, 'r') as f:
7     total_images = f['images'].shape[0]
8     indices = np.random.choice(total_images, num_samples_to_check, replace=False)
9     indices.sort() # Sort for faster image reading from HDF5 file
10
11    # Load sample images and check pixel values
12    pixel_values = []
```
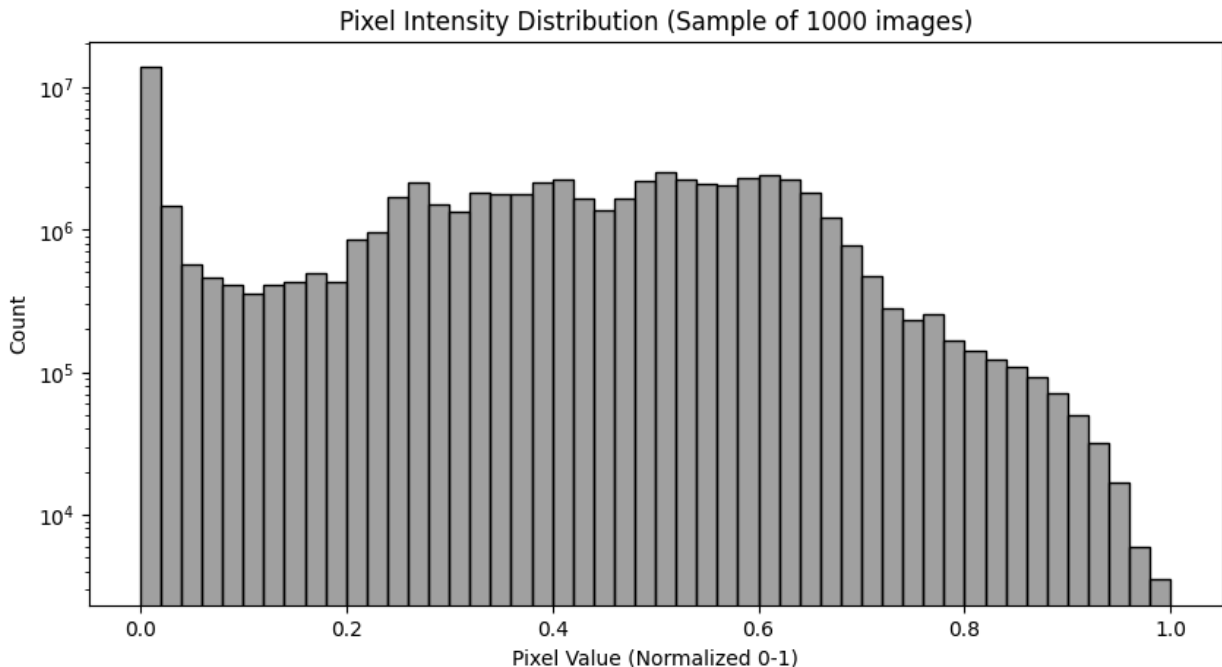
```
13      print(f'Loading {num_samples_to_check} random images...')
14      for idx in tqdm(indices):
15          img = f['images'][idx]
16          pixel_values.append(img.flatten())
17
18      print('Concatenating all pixel values...')
19      all_pixels = np.concatenate(pixel_values)
20
21 # Plot
22 print('Generating histogram plot...')
23 plt.figure(figsize=(10, 5))
24 sns.histplot(all_pixels, bins=50, color='gray', kde=False)
25 plt.title(f'Pixel Intensity Distribution (Sample of {num_samples_to_check} images)')
26 plt.xlabel('Pixel Value (Normalized 0-1)')
27 plt.ylabel('Count')
28 plt.yscale('log') # Log scale helps see the details in dark/bright areas
29 plt.show()
30
31 print(f'Min Pixel Value: {all_pixels.min():.4f}')
32 print(f'Max Pixel Value: {all_pixels.max():.4f}')
33 print(f'Mean Pixel Value: {all_pixels.mean():.4f}')
```

```
Loading 1000 random images...
  0%|          | 0/1000 [00:00<?, ?it/s]
Concatenating all pixel values...
Generating histogram plot...
```



Pixel Intensity Distribution (Sample of 1000 images)

```
Min Pixel Value: 0.0000
Max Pixel Value: 1.0000
Mean Pixel Value: 0.3466
```

```
 1 # Check patient id overlapping
 2
 3 # Load split IDs and Patient UIDs
 4 with h5py.File(data_path, 'r') as f:
 5     splits = f['split'][:] # 0=Train, 1=Val, 2=Test
 6     uids = f['StudyInstanceUID'][:]
 7
 8 # Get unique patients in each split
 9 train_patients = set(uids[splits == 0])
10 val_patients = set(uids[splits == 1])
11 test_patients = set(uids[splits == 2])
12
13 # Check for intersections
14 train_val_overlap = train_patients.intersection(val_patients)
15 train_test_overlap = train_patients.intersection(test_patients)
```

```
16 val_test_overlap = val_patients.intersection(test_patients)
17
18 print(f'Unique Patients in Train: {len(train_patients)}')
19 print(f'Unique Patients in Val:   {len(val_patients)}')
20 print(f'Unique Patients in Test:  {len(test_patients)}')
21 print("-" * 30)
22 print(f'Train-Val Overlap:  {len(train_val_overlap)} patients')
23 print(f'Train-Test Overlap: {len(train_test_overlap)} patients')
24 print(f'Val-Test Overlap:   {len(val_test_overlap)} patients')
25
26 if len(train_val_overlap) + len(train_test_overlap) + len(val_test_overlap) == 0:
27     print('\n✅ SUCCESS: No patient leakage detected between splits!')
28 else:
29     print('\n❌ WARNING: Patient leakage detected! Your splits are not independent.')
```

```
Unique Patients in Train: 1411
Unique Patients in Val:   302
Unique Patients in Test:  303
------------------------------
Train-Val Overlap:  0 patients
Train-Test Overlap: 0 patients
Val-Test Overlap:   0 patients

✅ SUCCESS: No patient leakage detected between splits!
```

```
1 # Verify image dimensions (should be 256x256)
2 with h5py.File(data_path, 'r') as f:
3   # Also check the dataset definition itself
4   dset = f['images']
5   print(f"Full Dataset Shape: {dset.shape}")
6   print(f"Full Dataset Chunks: {dset.chunks}")
```

```
Full Dataset Shape: (28812, 256, 256)
Full Dataset Chunks: (32, 256, 256)
```