# Dataloader Tutorial (Kaggle Notebook)

- To start a new notebook for model training, open a new Kaggle Notebook.
- Add the created private dataset by clicking "+ Add Input" on the righ-side pane. You should see the fracture_dataset_subset.h5 file there.

This notebook demonstrates how to go about loading the data from the private Kaggle dataset you had created in notebook 00. Here, we show you how to import the GitHub repository dataloader script into your notebook on Kaggle, how to view data samples from the preprocessed h5 file, and how to get the dataloaders based on the task at hand: 'classification' or 'detection'.

```python
1 # Import libraries
2 import h5py
3 import numpy as np
4 import os
5 import sys
6 import matplotlib.pyplot as plt
7 import matplotlib.patches as patches
8 import random
```

```python
1 # Import GitHub repo
2 # Repo files will appear in the /kaggle/working folder
3 !git clone https://github.com/apmalinsky/AAI-590-Capstone.git
4
5 scripts_path = '/kaggle/working/AAI-590-Capstone/scripts'
6
7 # Add the scripts folder to system path for importing
8 if scripts_path not in sys.path:
9     sys.path.append(scripts_path)
10    print(f'Added {scripts_path} to Python path.')
11 else:
12    print(f'{scripts_path} is already in the path.')
13
14 # Import data loader helper script
15 from dataloader import get_dataloaders
16 print('Dataloader script loaded!')
```

```
Cloning into 'AAI-590-Capstone'...
remote: Enumerating objects: 64, done.
remote: Counting objects: 100% (54/54), done.
remote: Compressing objects: 100% (45/45), done.
remote: Total 64 (delta 18), reused 26 (delta 6), pack-reused 10 (from 3)
Receiving objects: 100% (64/64), 2.34 MiB | 19.00 MiB/s, done.
Resolving deltas: 100% (18/18), done.
Added /kaggle/working/AAI-590-Capstone/scripts to Python path.
Dataloader script loaded!
```

```python
1 # View data file
2 data_path = '/kaggle/input/rsna-2022-hdf5-subset/fracture_dataset_subset.h5'
3
4 # Open the file
5 with h5py.File(data_path, 'r') as f:
6     # View key list, essentially our 'columns'
7     print(f'Keys: {list(f.keys())}')
8
9     # View column shapes
10    print('\n--- Details ---')
11    for key in f.keys():
12        dataset = f[key]
13        print(f'\nKey: {key}')
14        print(f'Shape: {dataset.shape}')
15        print(f'Dtype: {dataset.dtype}')
```

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

--- Details ---

Key: SliceNumber
Shape: (28868,)
Dtype: object

Key: StudyInstanceUID
Shape: (28868,)
Dtype: object

Key: bboxes
Shape: (28868, 10, 4)
Dtype: float32

Key: images
Shape: (28868, 256, 256)
Dtype: float32

Key: labels
Shape: (28868,)
Dtype: int8

Key: split
Shape: (28868,)
Dtype: int8
```

```python
 1 # View the first sample data
 2 with h5py.File(data_path, 'r') as f:
 3     slice_number = f['SliceNumber'][0]
 4     study_instance_uid = f['StudyInstanceUID'][0]
 5     bbox = f['bboxes'][0]
 6     image = f['images'][0]
 7     label = f['labels'][0]
 8     split = f['split'][0]
 9
10 print(f'SliceNumber: {slice_number}')
11 print(f'StudyInstanceUID: {study_instance_uid}')
12 print(f'Boxes: {bbox} (-1 is padding, which means no box)')
13 print(f'Image shape: {image.shape}')
14 print(f'Label: {label} (0=no fracture, 1=fracture)')
15 print(f'Split: {split} (0=train, 1=validation, 2=test)')
```

```
SliceNumber: b'59'
StudyInstanceUID: b'1.2.826.0.1.3680043.14345'
Boxes: [[-1. -1. -1. -1.]
 [-1. -1. -1. -1.]
 [-1. -1. -1. -1.]
 [-1. -1. -1. -1.]
 [-1. -1. -1. -1.]
 [-1. -1. -1. -1.]
 [-1. -1. -1. -1.]
 [-1. -1. -1. -1.]
 [-1. -1. -1. -1.]
 [-1. -1. -1. -1.]] (-1 is padding, which means no box)
Image shape: (256, 256)
Label: 0 (0=no fracture, 1=fracture)
Split: 0 (0=train, 1=validation, 2=test)
```

```python
 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()
```
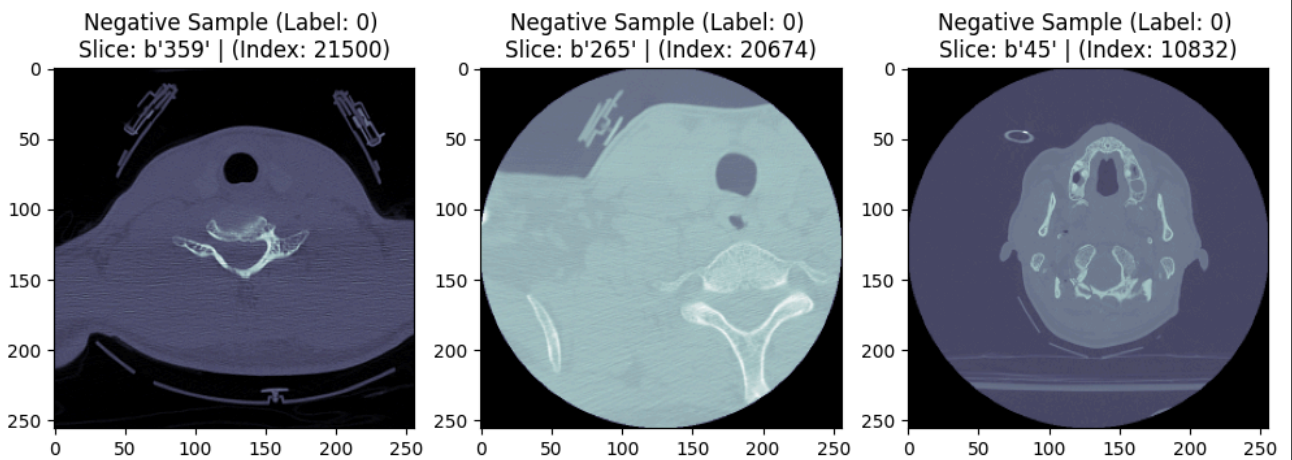
```
1 # View 3 sample negative images
2 visualize_samples(3, 0)
```
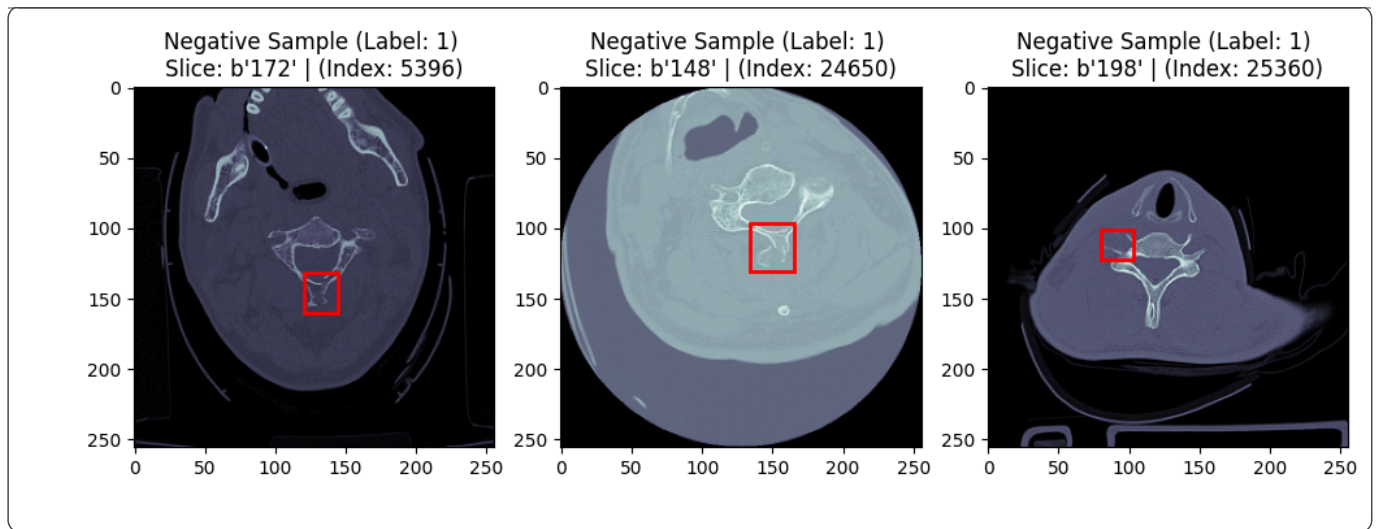


Negative Sample (Label: 0) Slice: b'359' | (Index: 21500)
Negative Sample (Label: 0) Slice: b'265' | (Index: 20674)
Negative Sample (Label: 0) Slice: b'45' | (Index: 10832)

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

Negative Sample (Label: 1) Slice: b'172' | (Index: 5396)
Negative Sample (Label: 1) Slice: b'148' | (Index: 24650)
Negative Sample (Label: 1) Slice: b'198' | (Index: 25360)

## Classification Dataloaders

- Returns: (images, labels)
- Images: A single tensor: [BatchSize, 1, 256, 256]
- Labels: A single tensor: [BatchSize] (with values 0 or 1)

```
1 # Get data loaders for classification
2 batch_size = 32
3
4 train_loader_cls, val_loader_cls, test_loader_cls = get_dataloaders(
5     hdf5_path=data_path,
6     batch_size=batch_size,
7     task='classification'
8 )
```

```
Loading dataset for task: 'classification'...
Loading pre-defined train/val/test splits...
Total samples: 28868
  Training indices:   20207
  Validation indices: 4330
  Test indices:       4331
DataLoaders created successfully.
```

## Detection Dataloaders

- Returns: (images, targets)
- Images: A list of image tensors (length BatchSize)
- Targets: A list of target dictionaries (length BatchSize)

```
1 # Get data loaders for detection
2 batch_size = 8 # lower batch size since detection models need more memory
3
4 train_loader_det, val_loader_det, test_loader_det = get_dataloaders(
5     hdf5_path=data_path,
6     batch_size=batch_size,
7     task='detection'
8 )
```

```
Loading dataset for task: 'detection'...
Loading pre-defined train/val/test splits...
Total samples: 28868
  Training indices:   20207
  Validation indices: 4330
  Test indices:       4331
DataLoaders created successfully.
```

## Use the Dataloaders for Model Training and Evaluation

1. Get the dataloaders for your task (classification vs detection)
2. Define your model
3. Perform training loop

   - Load images in batches "for images, labels in tqdm(train_loader_cls, desc='Training Batch'):"