

## ✓ Dataloader Setup for Colab Notebooks

```
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
```

```
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:

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

100% 3.51G/3.51G [00:40<00:00, 262MB/s]

100% 3.51G/3.51G [00:40<00:00, 92.9MB/s]

```
1 # Unzip and move to the local Colab storage for training
2 zip_name = KAGGLE_DATASET_ID.split('/')[-1] + '.zip'
3 !unzip -q '{zip_name}'
4
5 # This is the path the training script will use
6 HDF5_FILE_PATH = '/content/fracture_dataset_subset.h5'
```

```
1 # Get dataloader script from GitHub
2 !wget https://raw.githubusercontent.com/apmalinsky/AAI-590-Capstone/refs/heads/main/scripts/dataloader.py
```

```
--2025-11-19 14:26:18-- https://raw.githubusercontent.com/apmalinsky/AAI-590-Capstone/refs/heads/main/scripts/data
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 185.199.108.133, 185.199.109.133, 185.199.110.13
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|185.199.108.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
```

Length: 5917 (5.8K) [text/plain]  
 Saving to: 'dataloader.py'

dataloader.py 100%[=====>] 5.78K --.-KB/s in 0s

2025-11-19 14:26:18 (75.2 MB/s) - 'dataloader.py' saved [5917/5917]

```
1 # Import dataloader script
2 from dataloader import get_dataloaders
```

```
1 # View data file
2 data_path = '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

```
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)

```

```

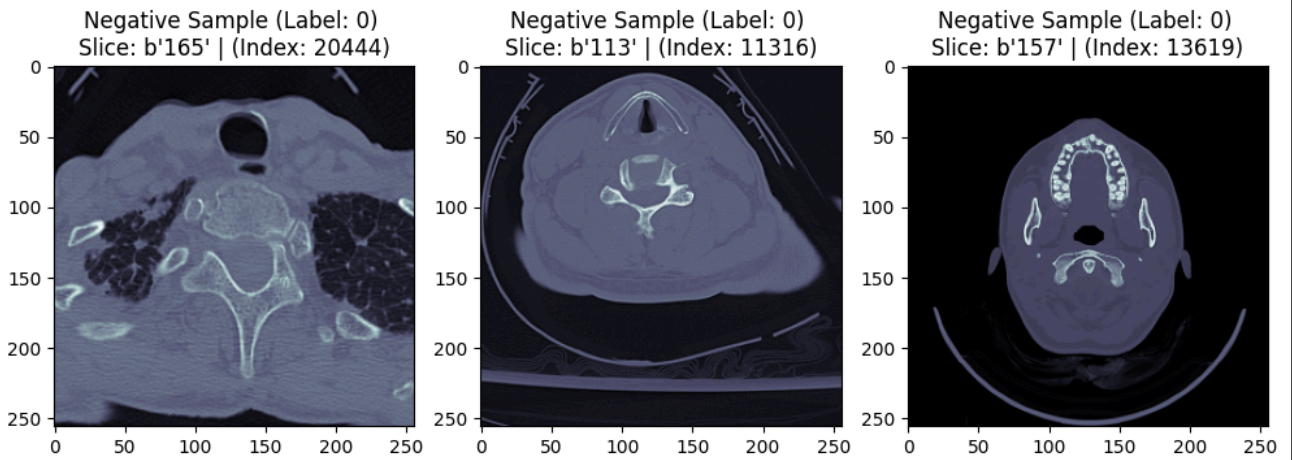
1 # Define visualizer function
2 def visualize_samples(num_samples, label):
3     with h5py.File(data_path, 'r') as hf:
4         # Load all labels and find the indices for the desired sample type
5         all_labels = hf['labels'][:]
6         indices_to_sample_from = np.where(all_labels == label)[0]
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()

```

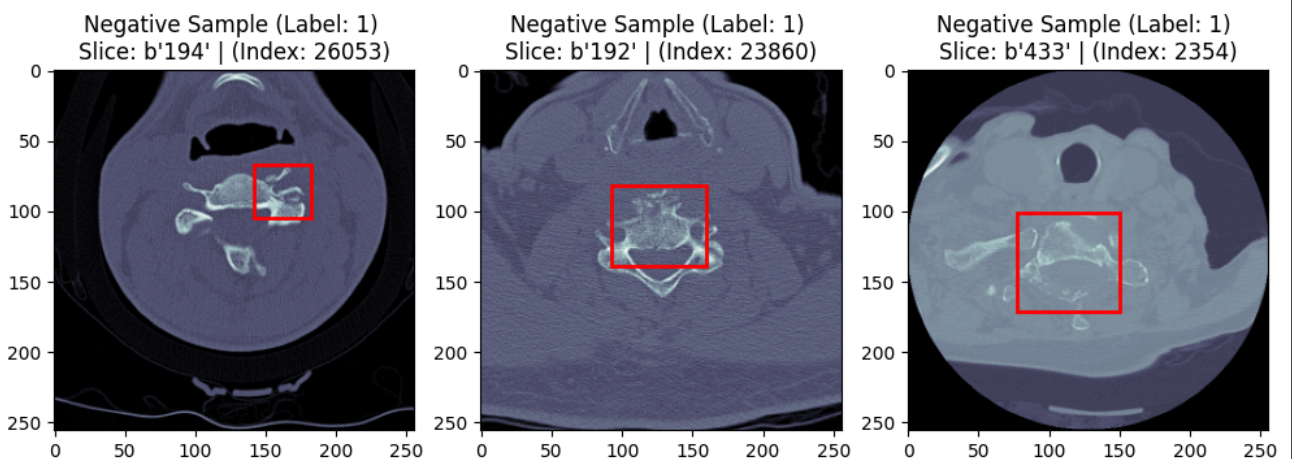
```

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

```



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



## 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'):"