

Create Your Own Private Subset Dataset for Cervical Spine Fracture Detection (Kaggle Notebook)

Before you get started, double-check you have the datasets loaded as Input (check right-hand side). You want to make sure you have these datasets selected:

- rsna-2022-spine-fracture-detection-metadata
- RSNA 2022 Cervical Spine Fracture Detection

```
1 # This will add the pydicom decompressors
2 !pip install pydicom pylibjpeg pylibjpeg-libjpeg
3 # Re-install a compatible version of numpy
4 !pip install numpy==1.26.4
```

!STOP AND RESTART SESSION!!

YOU ONLY NEED TO DO THIS ONCE

This makes sure the above dependencies are installed, and we're able to decompress images.

You can skip the first two cells after restarting.

```
1 # Import libraries
2 import pandas as pd
3 import numpy as np
4 import pydicom
5 import h5py
6 import cv2
7 import os
8 from tqdm.auto import tqdm
9 import warnings
10 import json
11 from sklearn.model_selection import train_test_split
12 from sklearn.model_selection import GroupShuffleSplit
13 import matplotlib.pyplot as plt
14 import seaborn as sns
15
16 random_state = 24
17
18 # Suppress the DeprecationWarning
19 warnings.filterwarnings("ignore", category=DeprecationWarning)
```

```
1 # Set paths to the Kaggle data
2 bbox_csv_path = '/kaggle/input/rsna-2022-cervical-spine-fracture-detection/train_bounding_boxes.csv'
3
4 # We're using the clean version from the metadata dataset
5 meta_csv_path = '/kaggle/input/rsna-2022-spine-fracture-detection-metadata/meta_train_clean.csv'
6
7 # This is the path to all 136GB of images
8 image_data_path = '/kaggle/input/rsna-2022-cervical-spine-fracture-detection/train_images/'
9
10 # Set the output path
11 hdf5_save_path = '/kaggle/working/fracture_dataset_subset.h5'
12
13 # Settings
14 image_size = 256
```

```
15 negative_ratio = 3 # for each positive, we add 3 negative samples
16 max_boxes = 10
```

```
1 # Load and view metadata
2 meta_df = pd.read_csv(meta_csv_path)
3 meta_df.head(3)
```

	StudyInstanceUID	Slice	ImageHeight	ImageWidth	SliceThickness	ImagePositionPatient_x	ImagePositionPatient_y	ImagePositionPatient_z
0	1.2.826.0.1.3680043.10001	1	512	512	0.625	-52.308	-100.0	-100.0
1	1.2.826.0.1.3680043.10001	2	512	512	0.625	-52.308	-100.0	-100.0
2	1.2.826.0.1.3680043.10001	3	512	512	0.625	-52.308	-100.0	-100.0

```
1 # Load and view bounding box data
2 bbox_df = pd.read_csv(bbox_csv_path)
3 bbox_df.head(3)
```

	StudyInstanceUID	x	y	width	height	slice_number
0	1.2.826.0.1.3680043.10051	219.27715	216.71419	17.30440	20.38517	133
1	1.2.826.0.1.3680043.10051	221.56460	216.71419	17.87844	25.24362	134
2	1.2.826.0.1.3680043.10051	216.82151	221.62546	27.00959	26.37454	135

```
1 # Load and merge metadata
2
3 # Rename slice columns to match
4 meta_df = meta_df.rename(columns={'Slice': 'SliceNumber'})
5 bbox_df = bbox_df.rename(columns={'slice_number': 'SliceNumber'})
6
7 # Group by slice and aggregate bboxes into a list [x, y, width, height]
8 print('Aggregating bounding boxes per slice...')
9 bbox_grouped = bbox_df.groupby(['StudyInstanceUID', 'SliceNumber'])[['x', 'y', 'width', 'height']].apply(
10     lambda x: x.values.tolist())
11 ).reset_index(name='bboxes')
12 print(f"Found {len(bbox_grouped)} slices with bounding boxes.")
13
14 # Merge all slices with the bounding box data
15 all_slices_df = pd.merge(
16     meta_df,
17     bbox_grouped,
18     on=['StudyInstanceUID', 'SliceNumber'],
19     how='left'
20 )
21
22 # Create final labels
23 # 'is_positive' is 1 if 'bboxes' is not NaN, 0 otherwise
24 all_slices_df['is_positive'] = all_slices_df['bboxes'].notna().astype(int)
25
26 # Fill NaN in 'bboxes' with an empty list for consistency
27 all_slices_df['bboxes'] = all_slices_df['bboxes'].apply(
28     lambda x: x if isinstance(x, list) else []
29 )
30
31 print('Metadata merge complete.')
```

Aggregating bounding boxes per slice...
Found 7217 slices with bounding boxes.
Metadata merge complete.

```
1 # Check image sizes
2 print(all_slices_df['ImageHeight'].value_counts())
3 print('\n',all_slices_df['ImageWidth'].value_counts())
```

```
1 ImageHeight
2 512    710819
3 768    782
4 Name: count, dtype: int64
```

```
1 ImageWidth
2 512    710574
3 768    782
4 519    245
5 Name: count, dtype: int64
```

```
1 # Filter for 512x512 images only
2 all_slices_df = all_slices_df[
3     (all_slices_df['ImageHeight'] == 512) &
4     (all_slices_df['ImageWidth'] == 512)
5 ]
6 print(all_slices_df['ImageHeight'].value_counts())
7 print('\n', all_slices_df['ImageWidth'].value_counts())
```

```
1 ImageHeight
2 512    710574
3 Name: count, dtype: int64
```

```
1 ImageWidth
2 512    710574
3 Name: count, dtype: int64
```

```
1 # Create patient-level splits
2
3 # Define Split Ratios
4 # 70% Train, 15% Validation, 15% Test
5 train_size = 0.70
6 val_size = 0.15
7 test_size = 0.15
8
9 # Prepare Patient-Level Data
10 # One row per patient for splitting
11 patient_df = all_slices_df.groupby('StudyInstanceUID').agg({
12     'is_positive': 'max', # If patient has any fracture slice, they are positive
13     'SliceNumber': 'count' # Count slices per patient
14 }).reset_index()
15
16 print(f'Total Patients: {len(patient_df)}')
17 print(f'Positive Patients: {patient_df['is_positive'].sum()}')
```

```
Total Patients: 2016
Positive Patients: 234
```

```
1 # First split: train / (val + test)
2
3 # We use GroupShuffleSplit to split by patient id while keeping class balance
4 splitter = GroupShuffleSplit(test_size=(val_size + test_size), n_splits=1, random_state=42)
5 train_idxs, temp_idxs = next(splitter.split(patient_df, y=patient_df['is_positive']), groups=patient_df['StudyI
6 train_patients = patient_df.iloc[train_idxs]
7 temp_patients = patient_df.iloc[temp_idxs]
8
9 # Second Split: val and test
10 # Split the remaining 30% into two equal halves (15% Val, 15% Test)
11 splitter_2 = GroupShuffleSplit(test_size=0.5, n_splits=1, random_state=42)
12 val_idxs, test_idxs = next(splitter_2.split(temp_patients, y=temp_patients['is_positive']), groups=temp_patient
13
14 val_patients = temp_patients.iloc[val_idxs]
15 test_patients = temp_patients.iloc[test_idxs]
16
17 print(f'Train Patient Count: {len(train_patients)})')
18 print(f'Validation Patient Count: {len(val_patients)})')
19 print(f'Test Patient Count: {len(test_patients)})')
```

```
Train Patient Count: 1411
Validation Patient Count: 302
```

Test Patient Count: 303

```

1 # Map splits back to slice dataframe based on patient id
2 # 0 = Train, 1 = Val, 2 = Test
3 train_uids = set(train_patients['StudyInstanceUID'])
4 val_uids = set(val_patients['StudyInstanceUID'])
5 test_uids = set(test_patients['StudyInstanceUID'])
6
7 def assign_split(uid):
8     if uid in train_uids: return 0
9     if uid in val_uids: return 1
10    if uid in test_uids: return 2
11
12 all_slices_df['split'] = all_slices_df['StudyInstanceUID'].apply(assign_split)
13 print('Split column created.')

```

Split column created.

```

1 # Create the Final Subset DataFrame
2
3 # Separate positives and negatives
4 positive_slices = all_slices_df[all_slices_df['is_positive'] == 1]
5 negative_slices = all_slices_df[all_slices_df['is_positive'] == 0]
6
7 # Sample negatives
8 negative_subset_list = []
9 for split_id in [0, 1, 2]:
10    split_negatives = negative_slices[negative_slices['split'] == split_id]
11    split_positives = positive_slices[positive_slices['split'] == split_id]
12
13    # Use our specified ratio
14    n_neg = int(len(split_positives) * negative_ratio)
15    n_neg = min(n_neg, len(split_negatives))
16    negative_subset_list.append(split_negatives.sample(n=n_neg, random_state=42))
17
18 negative_subset = pd.concat(negative_subset_list)
19
20 # Combine everything
21 final_subset_df = pd.concat([positive_slices, negative_subset]).sample(frac=1, random_state=42).reset_index()
22 final_subset_df.head()

```

	StudyInstanceUID	SliceNumber	ImageHeight	ImageWidth	SliceThickness	ImagePositionPatient_x	ImagePositionPatient_y	ImagePositionPatient_z
0	1.2.826.0.1.3680043.25600	155	512	512	0.625	-71.000000	-115.009300	-46.832031
1	1.2.826.0.1.3680043.24045	281	512	512	0.500	-64.000000	-79.921700	-71.000000
2	1.2.826.0.1.3680043.14348	352	512	512	0.600	-115.009300	-71.000000	-71.000000
3	1.2.826.0.1.3680043.12632	405	512	512	0.625	-71.000000	-71.000000	-71.000000
4	1.2.826.0.1.3680043.14087	364	512	512	0.500	-71.000000	-71.000000	-71.000000

```

1 # Verification
2 print(f"Train Slices: {len(final_subset_df[final_subset_df['split']==0])}")
3 print(f"Val Slices: {len(final_subset_df[final_subset_df['split']==1])}")
4 print(f"Test Slices: {len(final_subset_df[final_subset_df['split']==2])}")
5
6 # Check for leakage
7 train_uids = set(final_subset_df[final_subset_df['split']==0]['StudyInstanceUID'])
8 val_uids = set(final_subset_df[final_subset_df['split']==1]['StudyInstanceUID'])
9 test_uids = set(final_subset_df[final_subset_df['split']==2]['StudyInstanceUID'])
10
11 if len(train_uids.intersection(val_uids)) == 0 and len(train_uids.intersection(test_uids)) == 0:
12     print('✅ SUCCESS: No patient leakage detected!')
13 else:
14     print('❌ WARNING: Leakage detected!')

```

```
Train Slices: 20744
Val Slices: 4180
Test Slices: 3888
✓ SUCCESS: No patient leakage detected!
```

```
1 # View boxes before box scaling
2 final_subset_df['bboxes']
```

```
0 []
1 []
2 []
3 []
4 []
...
28807 []
28808 [[257.0, 160.0, 126.0, 95.0]]
28809 [[235.76963, 221.0, 84.23037, 83.94740999999999]]
28810 []
28811 []
Name: bboxes, Length: 28812, dtype: object
```

```
1 # Scale bounding boxes from original 512x512 to 256x256
2 # We can divide each box coordinate by 2, or multiply by half
3 scale_factor = 0.5
4
5 def scale_bboxes(bbox_list, scale_factor):
6     if not bbox_list: # for negative sample
7         return []
8
9     scaled_boxes = []
10    for box in bbox_list:
11        # Scale each coordinate [x, y, w, h]
12        scaled_box = [coord * scale_factor for coord in box]
13        scaled_boxes.append(scaled_box)
14
15    return scaled_boxes
16
17 # Apply this function to the 'bboxes' column
18 final_subset_df['bboxes'] = final_subset_df['bboxes'].apply(
19     lambda x: scale_bboxes(x, scale_factor)
20 )
21
22 print('Bounding boxes successfully scaled down for 256x256 size.')
23 final_subset_df['bboxes'] # after box scaling
```

```
Bounding boxes successfully scaled down for 256x256 size.
0 []
1 []
2 []
3 []
4 []
...
28807 []
28808 [[128.5, 80.0, 63.0, 47.5]]
28809 [[117.884815, 110.5, 42.115185, 41.97370499999...]
28810 []
28811 []
Name: bboxes, Length: 28812, dtype: object
```

```
1 # Image preprocessing function
2 def load_and_process_dicom(uid, slice_num, target_size):
3     # Build the local file path
4     file_path = f'{image_data_path}/{uid}/{slice_num}.dcm'
5
6     ds = pydicom.dcmread(file_path)
7
8     # Get pixel array
9     img = ds.pixel_array.astype(np.float32)
10
11    # Perform min-max normalization (scales pixel intensity to range [0,1])
12    img_min = np.min(img)
```

```

13     img_max = np.max(img)
14     if img_max > img_min:
15         img = (img - img_min) / (img_max - img_min)
16     else:
17         img = np.zeros(img.shape) # Handle black images
18
19     # Resize
20     img = cv2.resize(img, (target_size, target_size), interpolation=cv2.INTER_LINEAR)
21
22     return img

```

```

1 ## This cell was made using Google Gemini AI ##
2
3 num_samples = len(final_subset_df)
4 print(f"Creating HDF5 file at {hdf5_save_path} with {num_samples} total samples...")
5
6 # We'll write in chunks of 32 images at a time
7 image_chunk = (32, image_size, image_size)
8 bbox_chunk = (32, max_boxes, 4)
9 text_chunk = (512,) # A good chunk size for 1D arrays
10
11 with h5py.File(hdf5_save_path, 'w') as hf:
12
13     # --- Create datasets ---
14     dset_images = hf.create_dataset('images', shape=(num_samples, image_size, image_size), dtype='f4', chunks=
15     dset_labels = hf.create_dataset('labels', shape=(num_samples,), dtype='i1', chunks=text_chunk)
16     dset_bboxes = hf.create_dataset('bboxes', shape=(num_samples, max_boxes, 4), dtype='f4', fillvalue=-1.0, c
17     dt_str = h5py.special_dtype(vlen=str)
18     dset_uid = hf.create_dataset('StudyInstanceUID', (num_samples,), dtype=dt_str, chunks=text_chunk)
19     dset_slice = hf.create_dataset('SliceNumber', (num_samples,), dtype=dt_str, chunks=text_chunk)
20     dset_split = hf.create_dataset('split', shape=(num_samples,), dtype='i1', chunks=text_chunk)
21
22     # --- Start the processing loop ---
23     print("Starting processing loop...")
24
25     for idx, row in tqdm(final_subset_df.iterrows(), total=num_samples, desc="Processing slices"):
26         try:
27             # 1. Load and process the image from the local disk
28             img = load_and_process_dicom(
29                 row['StudyInstanceUID'],
30                 row['SliceNumber'],
31                 image_size
32             )
33
34             # 2. Save data to HDF5 file
35             dset_images[idx] = img
36             dset_labels[idx] = row['is_positive']
37             dset_uid[idx] = row['StudyInstanceUID']
38             dset_slice[idx] = str(row['SliceNumber'])
39             dset_split[idx] = row['split']
40
41             # 3. Process and save bounding boxes
42             bboxes = row['bboxes']
43             num_boxes = min(len(bboxes), max_boxes)
44
45             if num_boxes > 0:
46                 dset_bboxes[idx, :num_boxes, :] = np.array(bboxes[:num_boxes])
47
48         except Exception as e:
49             # This will catch any corrupted files
50             print(f"\n[Warning] Failed to process slice {row['StudyInstanceUID']}/{row['SliceNumber']}: {e}")
51
52 print("\n--- HDF5 file creation complete! ---")
53 print(f"File saved to: {hdf5_save_path}")
54 !ls -lh /kaggle/working/

```

```
Creating HDF5 file at /kaggle/working/fracture_dataset_subset.h5 with 28812 total samples...
Starting processing loop...
Processing slices:  0%|          | 0/28812 [00:00<?, ?it/s]

--- HDF5 file creation complete! ---
File saved to: /kaggle/working/fracture_dataset_subset.h5
total 7.1G
-rw-r--r-- 1 root root 113 Nov 26 23:23 dataset-metadata.json
-rw-r--r-- 1 root root 7.1G Nov 27 00:24 fracture_dataset_subset.h5
```

```
1 # Count total positive and negative samples
2 with h5py.File(hdf5_save_path, 'r') as f:
3     labels = f['labels'][:]
4     total_negative = np.sum(labels == 0)
5     total_positive = np.sum(labels == 1)
6
7 print(f'Total Negative Samples: {total_negative}')
8 print(f'Total Positive Samples: {total_positive}')
```

```
Total Negative Samples: 21609
Total Positive Samples: 7203
```

⚡ Workaround: Creating a Private Dataset with `kaggle.json`

If you are experiencing issues with the Secret add-on feature in Kaggle Notebooks, you can use a **private Kaggle Dataset** to securely store and access your `kaggle.json` file.

Steps to Create the Dataset:

1. **Locate Your File:** Find the `kaggle.json` file you just downloaded on your computer.
2. **Navigate to Dataset Creation:** Go to the Kaggle website and click "Create" in the top navigation bar, then select "New Dataset".
3. **Upload the File:**
 - Click "Upload File".
 - Select and upload your `kaggle.json` file.
4. **Configure Dataset Settings:**
 - **Title:** Give your dataset a descriptive name, like `Kaggle API Key`.
 - **Visibility:** **CRITICALLY IMPORTANT:** Ensure the visibility is set to "**Private**". This keeps your API key secure.
5. **Create and Finalize:** Click the "Create" or "Publish" button to finalize the dataset creation.

Steps to Use the Dataset in Your Notebook:

1. **Open Your Notebook:** Open the Kaggle Notebook where you need to use the API key.
2. **Add the Data:** In the right-hand panel of the notebook editor, click on "+ Add data".
3. **Search and Select:** Search for the private dataset you just created (e.g., `Kaggle API Key`) under the "Your Datasets" tab and select it.
4. **Access the Key:** Once the dataset is added, your `kaggle.json` file will be accessible in the notebook's file system, typically under the path:

```
/kaggle/input/<your-dataset-name>/kaggle.json
```

You can then modify your script to load the credentials from this path instead of relying on the Secret add-on.

```
1 # Load API credentials from your private api key dataset
2 # If you do not have this private dataset setup, see project README for instructions
3 # Doing this to work around Secret add-ons not working
4
5 private_secret_dataset_name = 'DATASET_NAME' ### REPLACE WITH YOUR PRIVATE DATASET NAME ####
6
7 CREDENTIALS_PATH = f'/kaggle/input/{private_secret_dataset_name}/kaggle.json'
```

```
8
9 # Create the hidden .kaggle directory
10 !mkdir -p ~/.kaggle
11
12 # Copy your key from the private dataset to the correct location
13 !cp '{CREDENTIALS_PATH}' ~/.kaggle/kaggle.json
14
15 # Set the correct permissions for the file
16 !chmod 600 ~/.kaggle/kaggle.json
17
18 print('Kaggle API credentials are now in place.')
```

Kaggle API credentials are now in place.

```
1 # Create the private Kaggle dataset
2
3 kaggle_username = 'KAGGLE_USERNAME' ### Set KAGGLE_USERNAME ####
4
5
6 # Define your dataset metadata
7 dataset_metadata = {
8     "title": "RSNA 2022 HDF5 Subset",
9     "id": f"{kaggle_username}/rsna-2022-hdf5-subset",
10    "licenses": [
11        {
12            "name": "CC0-1.0"
13        }
14    ]
15 }
16
17 # Write the metadata file
18 with open('/kaggle/working/dataset-metadata.json', 'w') as f:
19     json.dump(dataset_metadata, f)
20
21 # Run the create command
22 # You may have to wait a few minutes for the dataset to fully load to the Kaggle API system
23 print('Starting dataset creation... This will upload 7.05 GB.')
24 !kaggle datasets create -p /kaggle/working/ -r zip
```

```
Starting dataset creation... This will upload 7.05 GB.
Starting upload for file fracture_dataset_subset.h5
100%|██████████| 7.05G/7.05G [02:53<00:00, 43.7MB/s]
Upload successful: fracture_dataset_subset.h5 (7GB)
Starting upload for file .virtual_documents.zip
100%|██████████| 22.0/22.0 [00:00<00:00, 46.9B/s]
Upload successful: .virtual_documents.zip (22B)
Your private Dataset is being created. Please check progress at https://www.kaggle.com/datasets/andymalinsky/rsna-2
```