

```
!pip install h5py tensorflow numpy matplotlib
```

```
Requirement already satisfied: h5py in /usr/local/lib/python3.10/dist-packages (3.12.1)
Requirement already satisfied: tensorflow in /usr/local/lib/python3.10/dist-packages (2.17.1)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (1.26.4)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (3.7.5)
Requirement already satisfied: absl-py>=1.0.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (1.4.0)
Requirement already satisfied: astunparse>=1.6.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (1.6.3)
Requirement already satisfied: flatbuffers>=24.3.25 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (24.3.25)
Requirement already satisfied: gast!=0.5.0,!0.5.1,!0.5.2,>=0.2.1 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (0.4.1)
Requirement already satisfied: google-pasta>=0.1.1 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (0.2.0)
Requirement already satisfied: libclang>=13.0.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (18.1.1)
Requirement already satisfied: ml-dtypes<0.5.0,>=0.3.1 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (0.4.1)
Requirement already satisfied: opt-einsum>=2.3.2 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (3.4.0)
Requirement already satisfied: packaging in /usr/local/lib/python3.10/dist-packages (from tensorflow) (24.2)
Requirement already satisfied: protobuf!=4.21.0,!4.21.1,!4.21.2,!4.21.3,!4.21.4,!4.21.5,<5.0.0dev,>=3.20.3 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (4.21.5)
Requirement already satisfied: requests<3,>=2.21.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (2.32.3)
Requirement already satisfied: setuptools in /usr/local/lib/python3.10/dist-packages (from tensorflow) (75.1.0)
Requirement already satisfied: six>=1.12.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (1.17.0)
Requirement already satisfied: termcolor>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (2.5.0)
Requirement already satisfied: typing-extensions>=3.6.6 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (4.12.2)
Requirement already satisfied: wrapt>=1.11.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (1.17.0)
Requirement already satisfied: grpcio<2.0,>=1.24.3 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (1.68.1)
Requirement already satisfied: tensorboard<2.18,>=2.17 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (2.17.1)
Requirement already satisfied: keras>=3.2.0 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (3.5.0)
Requirement already satisfied: tensorflow-io-gcs-filesystem>=0.23.1 in /usr/local/lib/python3.10/dist-packages (from tensorflow) (0.37.0)
Requirement already satisfied: mkl_fft in /usr/local/lib/python3.10/dist-packages (from numpy) (1.3.8)
Requirement already satisfied: mkl_random in /usr/local/lib/python3.10/dist-packages (from numpy) (1.2.4)
Requirement already satisfied: mkl_umath in /usr/local/lib/python3.10/dist-packages (from numpy) (0.1.1)
Requirement already satisfied: mkl in /usr/local/lib/python3.10/dist-packages (from numpy) (2025.0.1)
Requirement already satisfied: tbb4py in /usr/local/lib/python3.10/dist-packages (from numpy) (2022.0.0)
Requirement already satisfied: mkl-service in /usr/local/lib/python3.10/dist-packages (from numpy) (2.4.1)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (1.3.1)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (4.55.3)
Requirement already satisfied: kiwisolver>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (1.4.7)
Requirement already satisfied: pillow>=6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (11.0.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (3.2.0)
Requirement already satisfied: python-dateutil>=2.7 in /usr/local/lib/python3.10/dist-packages (from matplotlib) (2.9.0.post0)
Requirement already satisfied: wheel<1.0,>=0.23.0 in /usr/local/lib/python3.10/dist-packages (from astunparse>=1.6.0->tensorflow) (0.43.0)
Requirement already satisfied: rich in /usr/local/lib/python3.10/dist-packages (from keras>=3.2.0->tensorflow) (13.9.4)
Requirement already satisfied: namex in /usr/local/lib/python3.10/dist-packages (from keras>=3.2.0->tensorflow) (0.0.8)
Requirement already satisfied: optree in /usr/local/lib/python3.10/dist-packages (from keras>=3.2.0->tensorflow) (0.13.1)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.21.0->tensorflow) (3.4.0)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.21.0->tensorflow) (3.10.1)
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.21.0->tensorflow) (2.3.0)
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests<3,>=2.21.0->tensorflow) (2025.1.31)
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.18,>=2.17->tensorflow) (3.7.0)
Requirement already satisfied: tensorboard-data-server<0.8.0,>=0.7.0 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.18,>=2.17->tensorflow) (0.20.0)
Requirement already satisfied: werkzeug>=1.0.1 in /usr/local/lib/python3.10/dist-packages (from tensorboard<2.18,>=2.17->tensorflow) (3.1.0)
Requirement already satisfied: intel-openmp>=2024 in /usr/local/lib/python3.10/dist-packages (from mkl->numpy) (2024.2.0)
Requirement already satisfied: tbb==2022.* in /usr/local/lib/python3.10/dist-packages (from mkl->numpy) (2022.0.0)
Requirement already satisfied: tcmlib==1.* in /usr/local/lib/python3.10/dist-packages (from tbb==2022.*->mkl->numpy) (1.2.0)
Requirement already satisfied: intel-cmplr-lib-rt in /usr/local/lib/python3.10/dist-packages (from mkl_umath->numpy) (2024.2.0)
Requirement already satisfied: intel-cmplr-lib-ur==2024.2.0 in /usr/local/lib/python3.10/dist-packages (from intel-openmp>=2024->mkl->numpy) (2024.2.0)
Requirement already satisfied: MarkupSafe>=2.1.1 in /usr/local/lib/python3.10/dist-packages (from werkzeug>=1.0.1->tensorboard<2.18,>=2.17->tensorflow) (3.0.2)
Requirement already satisfied: markdown-it-py>=2.2.0 in /usr/local/lib/python3.10/dist-packages (from rich->keras>=3.2.0->tensorflow) (3.0.0)
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in /usr/local/lib/python3.10/dist-packages (from rich->keras>=3.2.0->tensorflow) (2.19.0)
Requirement already satisfied: mdurl~0.1 in /usr/local/lib/python3.10/dist-packages (from markdown-it-py>=2.2.0->rich->keras>=3.2.0->tensorflow) (0.1.2)
```

```
import h5py
```

```
file_path = '/kaggle/input/quark-gluon-ds/quark-gluon_data-set_n139306.hdf5'
```

```
# Step 1: Check if the file can be opened
```

```
try:
    with h5py.File(file_path, 'r') as f:
        print("File opened successfully.")
except Exception as e:
    print("Error opening file:", e)
```

```
# Step 2: List top-level keys
```

```
try:
    with h5py.File(file_path, 'r') as f:
        print("Top-level keys:", list(f.keys()))
except Exception as e:
    print("Error listing keys:", e)
```

```
# Step 3: Try accessing a specific dataset
```

```
try:
    with h5py.File(file_path, 'r') as f:
        if 'X_jets' in f:
            print("Shape of X_jets:", f['X_jets'].shape)
        else:
            print("X_jets not found in the file.")
```

```
except Exception as e:
    print("Error accessing dataset:", e)
```

File opened successfully.
 Top-level keys: ['X_jets', 'm0', 'pt', 'y']
 Shape of X_jets: (139306, 125, 125, 3)

```
import time
import threading
```

```
def prevent_disconnect():
    """
    ✓ Runs a background thread that prints a message every 5-10 minutes
    to prevent Kaggle from disconnecting.
    """
    interval = 60*10 # ✓ 5 minutes interval
    print("Preventing Kaggle auto-disconnect...")

    def keep_active():
        while True:
            print(" Keeping Kaggle active...")
            time.sleep(interval)

    # ✓ Run the keep-alive thread
    thread = threading.Thread(target=keep_active)
    thread.daemon = True # ✓ Stops the thread when the script stops
    thread.start()
```

```
# ✓ Start the script
prevent_disconnect()
```

Preventing Kaggle auto-disconnect...
 Keeping Kaggle active...

```
!pip install torch_geometric
```

Requirement already satisfied: torch_geometric in /usr/local/lib/python3.10/dist-packages (2.6.1)
 Requirement already satisfied: aiohttp in /usr/local/lib/python3.10/dist-packages (from torch_geometric) (3.11.12)
 Requirement already satisfied: fsspec in /usr/local/lib/python3.10/dist-packages (from torch_geometric) (2024.12.0)
 Requirement already satisfied: Jinja2 in /usr/local/lib/python3.10/dist-packages (from torch_geometric) (3.1.4)
 Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from torch_geometric) (1.26.4)
 Requirement already satisfied: psutil>=5.8.0 in /usr/local/lib/python3.10/dist-packages (from torch_geometric) (5.9.5)
 Requirement already satisfied: pyparsing in /usr/local/lib/python3.10/dist-packages (from torch_geometric) (3.2.0)
 Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from torch_geometric) (2.32.3)
 Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from torch_geometric) (4.67.1)
 Requirement already satisfied: aiohappyeyeballs>=2.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->torch_geometric) (2.4.4)
 Requirement already satisfied: aiosignal>=1.1.2 in /usr/local/lib/python3.10/dist-packages (from aiohttp->torch_geometric) (1.3.2)
 Requirement already satisfied: async-timeout<6.0,>=4.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->torch_geometric) (5.0.1)
 Requirement already satisfied: attrs>=17.3.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->torch_geometric) (25.1.0)
 Requirement already satisfied: frozenlist>=1.1.1 in /usr/local/lib/python3.10/dist-packages (from aiohttp->torch_geometric) (1.5.0)
 Requirement already satisfied: multidict<7.0,>=4.5 in /usr/local/lib/python3.10/dist-packages (from aiohttp->torch_geometric) (6.1.0)
 Requirement already satisfied: propcache>=0.2.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->torch_geometric) (0.2.1)
 Requirement already satisfied: yarl<2.0,>=1.17.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->torch_geometric) (1.18.3)
 Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from Jinja2->torch_geometric) (3.0.2)
 Requirement already satisfied: mkl_fft in /usr/local/lib/python3.10/dist-packages (from numpy->torch_geometric) (1.3.8)
 Requirement already satisfied: mkl_random in /usr/local/lib/python3.10/dist-packages (from numpy->torch_geometric) (1.2.4)
 Requirement already satisfied: mkl_umath in /usr/local/lib/python3.10/dist-packages (from numpy->torch_geometric) (0.1.1)
 Requirement already satisfied: mkl in /usr/local/lib/python3.10/dist-packages (from numpy->torch_geometric) (2025.0.1)
 Requirement already satisfied: tb4py in /usr/local/lib/python3.10/dist-packages (from numpy->torch_geometric) (2022.0.0)
 Requirement already satisfied: mkl-service in /usr/local/lib/python3.10/dist-packages (from numpy->torch_geometric) (2.4.1)
 Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->torch_geometric) (3.4.0)
 Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->torch_geometric) (3.10)
 Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->torch_geometric) (2.3.0)
 Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->torch_geometric) (2025.11.12)
 Requirement already satisfied: typing-extensions>=4.1.0 in /usr/local/lib/python3.10/dist-packages (from multidict<7.0,>=4.5->aiohttp) (4.13.2)
 Requirement already satisfied: intel-openmp>=2024 in /usr/local/lib/python3.10/dist-packages (from mkl->numpy->torch_geometric) (2024.2.0)
 Requirement already satisfied: tb==2022.* in /usr/local/lib/python3.10/dist-packages (from mkl->numpy->torch_geometric) (2022.0.0)
 Requirement already satisfied: tcmlib==1.* in /usr/local/lib/python3.10/dist-packages (from tb==2022.*->mkl->numpy->torch_geometric) (1.0.0)
 Requirement already satisfied: intel-cmplr-lib-rt in /usr/local/lib/python3.10/dist-packages (from mkl->numpy->torch_geometric) (2024.2.0)
 Requirement already satisfied: intel-cmplr-lib-ur==2024.2.0 in /usr/local/lib/python3.10/dist-packages (from intel-openmp>=2024->mkl)

```
!pip install tqdm scipy
```

Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (4.67.1)
 Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (1.13.1)
 Requirement already satisfied: numpy<2.3,>=1.22.4 in /usr/local/lib/python3.10/dist-packages (from scipy) (1.26.4)
 Requirement already satisfied: mkl_fft in /usr/local/lib/python3.10/dist-packages (from numpy<2.3,>=1.22.4->scipy) (1.3.8)
 Requirement already satisfied: mkl_random in /usr/local/lib/python3.10/dist-packages (from numpy<2.3,>=1.22.4->scipy) (1.2.4)
 Requirement already satisfied: mkl_umath in /usr/local/lib/python3.10/dist-packages (from numpy<2.3,>=1.22.4->scipy) (0.1.1)
 Requirement already satisfied: mkl in /usr/local/lib/python3.10/dist-packages (from numpy<2.3,>=1.22.4->scipy) (2025.0.1)
 Requirement already satisfied: tb4py in /usr/local/lib/python3.10/dist-packages (from numpy<2.3,>=1.22.4->scipy) (2022.0.0)
 Requirement already satisfied: mkl-service in /usr/local/lib/python3.10/dist-packages (from numpy<2.3,>=1.22.4->scipy) (2.4.1)
 Requirement already satisfied: intel-openmp>=2024 in /usr/local/lib/python3.10/dist-packages (from mkl->numpy<2.3,>=1.22.4->scipy)

Requirement already satisfied: tbb==2022.* in /usr/local/lib/python3.10/dist-packages (from mkl->numpy<2.3,>=1.22.4->scipy) (2022.0)

Requirement already satisfied: tcmlib==1.* in /usr/local/lib/python3.10/dist-packages (from tbb==2022.*->mkl->numpy<2.3,>=1.22.4->scipy) (1.0.1)

Requirement already satisfied: intel-cmplr-lib-rt in /usr/local/lib/python3.10/dist-packages (from mkl_umath->numpy<2.3,>=1.22.4->scipy) (2024.2.0)

Requirement already satisfied: intel-cmplr-lib-ur==2024.2.0 in /usr/local/lib/python3.10/dist-packages (from intel-openmp==2024->mkl) (2024.2.0)

```
import os
import numpy as np
import torch
from torch_geometric.data import Data
from scipy.spatial import KDTree
from tqdm import tqdm
import pandas as pd
from sklearn.preprocessing import MinMaxScaler

# ✓ Paths
point_cloud_dir = '/kaggle/working/PointClouds/'
graph_dir = '/kaggle/working/Graphs/'
os.makedirs(graph_dir, exist_ok=True)

# ✓ Parameters
K = 5 # Number of neighbors for KNN graph
RADIUS = 5.0 # Radius for connectivity

# ✓ Paths
batch_dir = "/kaggle/input/dataset-zip" # Path where batches are stored
output_dir = "/kaggle/working/PointClouds/" # Save point cloud dataset

# ✓ Ensure the output directory exists
os.makedirs(output_dir, exist_ok=True)

# ✓ Normalization Scaler
scaler = MinMaxScaler()

def convert_to_point_cloud(batch_file, output_file):
    """
    Converts a batch of jet images into a point cloud representation.
    Args:
        batch_file (str): Path to the batch .npz file.
        output_file (str): Path to save the point cloud dataset.
    """
    # ✓ Load the batch
    X_batch = np.load(batch_file) # Shape: (batch_size, 125, 125, 3)

    point_clouds = []

    for img_idx in range(X_batch.shape[0]):
        img = X_batch[img_idx] # Shape: (125, 125, 3)

        # ✓ Extract non-zero pixels
        non_zero_mask = np.any(img != 0, axis=-1) # True for non-zero pixels

        # ✓ Extract pixel coordinates (x, y)
        y_coords, x_coords = np.where(non_zero_mask)

        # ✓ Extract intensity features (ECAL, HCAL, Tracks)
        ecal = img[non_zero_mask][:, 0]
        hcal = img[non_zero_mask][:, 1]
        tracks = img[non_zero_mask][:, 2]

        # ✓ Normalize intensity features (0-1 range)
        features = np.column_stack((ecal, hcal, tracks))
        features_normalized = scaler.fit_transform(features)

        # ✓ Collect point cloud data
        for i in range(len(x_coords)):
            point_clouds.append([
                img_idx, # Event ID (image index)
                x_coords[i], # X-coordinate
                y_coords[i], # Y-coordinate
                features_normalized[i, 0], # ECAL normalized
                features_normalized[i, 1], # HCAL normalized
                features_normalized[i, 2] # Tracks normalized
            ])

    # ✓ Save as CSV for easy loading later
    df = pd.DataFrame(point_clouds, columns=["event_id", "x", "y", "ecal", "hcal", "tracks"])
    df.to_csv(output_file, index=False)
    print(f"✓ Point cloud saved: {output_file}")

# ✓ Iterate through all batches and convert them to point clouds
batch_files = sorted([f for f in os.listdir(batch_dir) if f.endswith('.npz')])
```

```

for i, batch_file in enumerate(tqdm(batch_files, desc="Converting Batches")):
    batch_path = os.path.join(batch_dir, batch_file)
    output_file = os.path.join(output_dir, f"point_cloud_{i+1}.csv")

    # ✅ Convert the batch to point cloud format
    convert_to_point_cloud(batch_path, output_file)

print(f"✅ All batches converted to point cloud format and saved in {output_dir}")

import os
import torch
import pandas as pd
import numpy as np
from scipy.spatial import cKDTree
from torch_geometric.data import Data
import h5py

# -----
# ✅ Constants
# -----
HDF5_FILE_PATH = '/kaggle/input/quark-gluon-ds/quark-gluon_data-set_n139306.hdf5'
BATCH_SIZE = 2000

# -----
# ✅ kNN Graph Construction Function
# -----
def knn_graph(x, k=5):
    """
    Constructs a kNN graph from node features.

    Args:
        x (torch.Tensor): Node features (size: [num_nodes, num_features]).
        k (int): Number of neighbors.

    Returns:
        edge_index (torch.Tensor): Tensor of shape [2, num_edges].
    """
    x_np = x.cpu().numpy()
    tree = cKDTree(x_np)

    _, neighbors = tree.query(x_np, k=k + 1) # k+1 to include self
    edge_index = []

    for i, n in enumerate(neighbors):
        for j in n[1:]: # Skip self-loop
            edge_index.append([i, j])
            edge_index.append([j, i])

    edge_index = torch.tensor(edge_index, dtype=torch.long).T
    return edge_index

# -----
# ✅ Function to Extract Labels in Batches
# -----
def extract_labels_in_batches(hdf5_file, batch_size=2000):
    """
    Extracts labels from the HDF5 file in batches.

    Args:
        hdf5_file (str): Path to the HDF5 file.
        batch_size (int): Size of each batch.

    Returns:
        List[torch.Tensor]: List of label tensors.
    """
    labels = []

    with h5py.File(hdf5_file, 'r') as f:
        y_data = f['y'] # Extract labels
        total_samples = y_data.shape[0]

        print(f"✅ Total Samples: {total_samples}")

        for i in range(0, total_samples, batch_size):
            batch_labels = torch.tensor(y_data[i:i + batch_size], dtype=torch.long)
            labels.append(batch_labels)

        print(f"🔥 Extracted batch {i // batch_size + 1} / {total_samples // batch_size + 1}")

    return torch.cat(labels, dim=0) # Combine all batches into one tensor

```

```

# -----
# ✅ Load Point Cloud CSVs and Create Graphs with Actual Labels
# -----
def load_and_cache_graphs_with_labels(data_dirs, labels, k=5, cache_dir='/kaggle/working/Cached_Graphs'):
    """
    Loads point cloud CSV files, builds graph objects with actual labels, and caches them.

    Args:
        data_dirs (list of str): List of directories containing CSVs.
        labels (torch.Tensor): Tensor of labels extracted from the HDF5 file.
        k (int): Number of neighbors for kNN graph.
        cache_dir (str): Directory to save cached graph files.

    Returns:
        List[Data]: List of PyG Data objects with actual labels.
    """
    os.makedirs(cache_dir, exist_ok=True)
    graph_list = []
    label_idx = 0 # Keep track of label assignment

    for data_dir in data_dirs:
        if not os.path.exists(data_dir):
            print(f"⚠️ Directory {data_dir} does not exist, skipping.")
            continue

        for file in sorted(os.listdir(data_dir)): # Ensure file order matches label order
            if file.endswith('.csv'):
                file_path = os.path.join(data_dir, file)
                cache_path = os.path.join(cache_dir, file + ".pt")

                # ✅ Check if already cached
                if os.path.exists(cache_path):
                    print(f"✅ Loading cached graph: {file}")
                    graph = torch.load(cache_path)
                else:
                    print(f"🔥 Processing and caching: {file}")
                    df = pd.read_csv(file_path)

                    # ✅ Extract features: [x, y, ecal, hcal, tracks]
                    x_tensor = torch.tensor(df[['x', 'y', 'ecal', 'hcal', 'tracks']].values, dtype=torch.float)

                    # ✅ Construct kNN graph
                    edge_index = knn_graph(x_tensor, k=k)

                    # ✅ Assign correct label
                    if label_idx < len(labels):
                        label = labels[label_idx].unsqueeze(0)
                        label_idx += 1
                    else:
                        label = torch.tensor([0], dtype=torch.long) # Fallback label

                    # ✅ Create PyG Data object
                    graph = Data(x=x_tensor, edge_index=edge_index, y=label)

                    # ✅ Cache the graph
                    torch.save(graph, cache_path)

                graph_list.append(graph)

        print(f"✅ Loaded {len(graph_list)} graphs with actual labels.")
    return graph_list

# -----
# Specify Directories and Load Graphs
# -----

if __name__ == "__main__":

    point_cloud_dir1 = '/kaggle/input/point-clouds'
    point_cloud_dir2 = '/kaggle/input/extra-pc-ds'
    data_dirs = [point_cloud_dir1, point_cloud_dir2]
    # ✅ Step 1: Extract labels from HDF5
    labels = extract_labels_in_batches(HDF5_FILE_PATH, batch_size=BATCH_SIZE)

    # ✅ Step 2: Load CSVs and create graphs with actual labels
    graphs = load_and_cache_graphs_with_labels(data_dirs, labels, k=5, cache_dir='/kaggle/working/Cached_Graphs')

    # # -----
    # # Display Basic Info of Loaded Graphs
    # # -----
    # if len(graph_data) > 0:

```

```
#     print(f"🔥 Total Graphs Loaded: {len(graph_data)}")
#     print(f"🔥 Example Graph Node Shape: {graph_data[0].x.shape}")
#     print(f"🔥 Example Graph Edge Index Shape: {graph_data[0].edge_index.shape}")
#     print(f"🔥 Example Graph Label: {graph_data[0].y.item()}")
# else:
#     print("🚩 No graphs loaded.")
```



✓ Total Samples: 139306
🔥 Extracted batch 1 / 70
🔥 Extracted batch 2 / 70
🔥 Extracted batch 3 / 70
🔥 Extracted batch 4 / 70
🔥 Extracted batch 5 / 70
🔥 Extracted batch 6 / 70
🔥 Extracted batch 7 / 70
🔥 Extracted batch 8 / 70
🔥 Extracted batch 9 / 70
🔥 Extracted batch 10 / 70
🔥 Extracted batch 11 / 70
🔥 Extracted batch 12 / 70
🔥 Extracted batch 13 / 70
🔥 Extracted batch 14 / 70
🔥 Extracted batch 15 / 70
🔥 Extracted batch 16 / 70
🔥 Extracted batch 17 / 70
🔥 Extracted batch 18 / 70
🔥 Extracted batch 19 / 70
🔥 Extracted batch 20 / 70
🔥 Extracted batch 21 / 70
🔥 Extracted batch 22 / 70
🔥 Extracted batch 23 / 70
🔥 Extracted batch 24 / 70
🔥 Extracted batch 25 / 70
🔥 Extracted batch 26 / 70
🔥 Extracted batch 27 / 70
🔥 Extracted batch 28 / 70
🔥 Extracted batch 29 / 70
🔥 Extracted batch 30 / 70
🔥 Extracted batch 31 / 70
🔥 Extracted batch 32 / 70
🔥 Extracted batch 33 / 70
🔥 Extracted batch 34 / 70
🔥 Extracted batch 35 / 70
🔥 Extracted batch 36 / 70
🔥 Extracted batch 37 / 70
🔥 Extracted batch 38 / 70
🔥 Extracted batch 39 / 70
🔥 Extracted batch 40 / 70
🔥 Extracted batch 41 / 70
🔥 Extracted batch 42 / 70
🔥 Extracted batch 43 / 70
🔥 Extracted batch 44 / 70
🔥 Extracted batch 45 / 70
🔥 Extracted batch 46 / 70
🔥 Extracted batch 47 / 70
🔥 Extracted batch 48 / 70
🔥 Extracted batch 49 / 70
🔥 Extracted batch 50 / 70
🔥 Extracted batch 51 / 70
🔥 Extracted batch 52 / 70
🔥 Extracted batch 53 / 70
🔥 Extracted batch 54 / 70
🔥 Extracted batch 55 / 70
🔥 Extracted batch 56 / 70
🔥 Extracted batch 57 / 70
🔥 Extracted batch 58 / 70
🔥 Extracted batch 59 / 70
🔥 Extracted batch 60 / 70
🔥 Extracted batch 61 / 70
🔥 Extracted batch 62 / 70
🔥 Extracted batch 63 / 70
🔥 Extracted batch 64 / 70
🔥 Extracted batch 65 / 70
🔥 Extracted batch 66 / 70
🔥 Extracted batch 67 / 70
🔥 Extracted batch 68 / 70
🔥 Extracted batch 69 / 70
🔥 Extracted batch 70 / 70
🔥 Processing and caching: point_cloud_1.csv
🔥 Processing and caching: point_cloud_10.csv
🔥 Processing and caching: point_cloud_11.csv
🔥 Processing and caching: point_cloud_12.csv
🔥 Processing and caching: point_cloud_14.csv
🔥 Processing and caching: point_cloud_15.csv
🔥 Processing and caching: point_cloud_16.csv
🔥 Processing and caching: point_cloud_17.csv
🔥 Processing and caching: point_cloud_18.csv
🔥 Processing and caching: point_cloud_19.csv
🔥 Processing and caching: point_cloud_2.csv
🔥 Processing and caching: point_cloud_20.csv
🔥 Processing and caching: point_cloud_21.csv
🔥 Processing and caching: point_cloud_22.csv
🔥 Processing and caching: point_cloud_23.csv
🔥 Processing and caching: point_cloud_24.csv
🔥 Processing and caching: point_cloud_25.csv
🔥 Processing and caching: point_cloud_26.csv
🔥 Processing and caching: point_cloud_27.csv

🟡 Processing and caching: point_cloud_28.csv
 Keeping Kaggle active...
 🟡 Processing and caching: point_cloud_29.csv
 🟡 Processing and caching: point_cloud_3.csv
 🟡 Processing and caching: point_cloud_30.csv
 🟡 Processing and caching: point_cloud_31.csv
 🟡 Processing and caching: point_cloud_32.csv
 🟡 Processing and caching: point_cloud_33.csv
 🟡 Processing and caching: point_cloud_34.csv
 🟡 Processing and caching: point_cloud_35.csv
 🟡 Processing and caching: point_cloud_36.csv
 🟡 Processing and caching: point_cloud_37.csv
 🟡 Processing and caching: point_cloud_38.csv
 🟡 Processing and caching: point_cloud_39.csv
 🟡 Processing and caching: point_cloud_4.csv
 🟡 Processing and caching: point_cloud_40.csv
 🟡 Processing and caching: point_cloud_41.csv
 🟡 Processing and caching: point_cloud_42.csv
 🟡 Processing and caching: point_cloud_43.csv
 🟡 Processing and caching: point_cloud_44.csv
 🟡 Processing and caching: point_cloud_45.csv
 🟡 Processing and caching: point_cloud_46.csv
 🟡 Processing and caching: point_cloud_47.csv
 Keeping Kaggle active...
 🟡 Processing and caching: point_cloud_48.csv
 🟡 Processing and caching: point_cloud_49.csv
 🟡 Processing and caching: point_cloud_5.csv
 🟡 Processing and caching: point_cloud_50.csv
 🟡 Processing and caching: point_cloud_6.csv
 🟡 Processing and caching: point_cloud_7.csv
 🟡 Processing and caching: point_cloud_8.csv
 🟡 Processing and caching: point_cloud_9.csv
 🟡 Processing and caching: point_cloud_13.csv
 ✅ Loaded 50 graphs with actual labels.

```

-----
NameError                                Traceback (most recent call last)
<ipython-input-8-48596739762c> in <cell line: 20>()
    18 # Display Basic Info of Loaded Graphs
    19 # -----
--> 20 if len(graph_data) > 0:
    21     print(f"🟡 Total Graphs Loaded: {len(graph_data)}")
    22     print(f"🟡 Example Graph Node Shape: {graph_data[0].x.shape}")
  
```

NameError: name 'graph_data' is not defined

Keeping Kaggle active...


```
import os
```

```
cache_folder = "/kaggle/input/cached-graphs"
files = os.listdir(cache_folder)
print("Number of files in Cached_Graphs:", len(files))
print("Files:", files)
```

Number of files in Cached_Graphs: 50
Files: ['point_cloud_24.csv.pt', 'point_cloud_48.csv.pt', 'point_cloud_41.csv.pt', 'point_cloud_45.csv.pt', 'point_cloud_17.csv.pt',

```
import torch
import torch.nn.functional as F
from torch.nn import Linear, BatchNorm1d
from torch_geometric.data import DataLoader, Batch
import os
import time
import matplotlib.pyplot as plt
import gc

# -----
# Graph Convolution Layer (Manual)
# -----
# ☒ Ensure consistent tensor types inside the GraphConv class
class GraphConv(torch.nn.Module):
    def __init__(self, in_channels, out_channels):
        super(GraphConv, self).__init__()
        self.lin = torch.nn.Linear(in_channels, out_channels)
        self.bn = torch.nn.BatchNorm1d(out_channels)

    def forward(self, x, edge_index, edge_weight=None):
        row, col = edge_index

        if edge_weight is None:
            edge_weight = torch.ones(row.size(0), dtype=torch.float32, device=x.device)

        # ☒ Explicitly cast tensors to float32 before index_add_
        x = x.float() # Ensure x is in float32
        edge_weight = edge_weight.float() # Ensure weights are in float32

        # ☒ Aggregate neighbor messages
        out = torch.zeros_like(x)
        out.index_add_(0, row, x[col] * edge_weight.view(-1, 1))

        out = self.bn(self.lin(out))
        return F.relu(out)

# -----
# DGCNN Model Definition
# -----
class DGCNN(torch.nn.Module):
    def __init__(self, num_features=5, hidden_dim=64):
        super(DGCNN, self).__init__()
        self.conv1 = GraphConv(num_features, hidden_dim)
        self.conv2 = GraphConv(hidden_dim, hidden_dim)
        # MLP layers for classification
        self.lin1 = Linear(hidden_dim, 128)
        self.lin2 = Linear(128, 64)
        self.lin3 = Linear(64, 2)

    def forward(self, data):
        # data is a Batch object from PyG
        x, edge_index, batch = data.x, data.edge_index, data.batch
        x = self.conv1(x, edge_index)
        x = self.conv2(x, edge_index)
        # Global pooling: using mean pooling for each graph
        x = torch_geometric.nn.global_mean_pool(x, batch)
        x = F.relu(self.lin1(x))
        x = F.relu(self.lin2(x))
        x = self.lin3(x)
        return F.log_softmax(x, dim=1)

import os
import torch
from torch_geometric.data import Data, DataLoader
from torch.serialization import add_safe_globals

# ☒ Allowlist the required PyG class
add_safe_globals([Data])
```

```

# ✅ Path to Cached Graphs
cache_folder = "/kaggle/input/cached-graphs"

# ✅ Lazy Loading Dataset
class GraphDataset(torch.utils.data.Dataset):
    def __init__(self, folder):
        self.folder = folder
        self.files = [f for f in os.listdir(folder) if f.endswith(".pt")]

    def __len__(self):
        return len(self.files)

    def __getitem__(self, idx):
        file_path = os.path.join(self.folder, self.files[idx])
        graph = torch.load(file_path, weights_only=False)
        return graph

# ✅ Load Graphs from Cached Files with `weights_only=False`
graph_data = []
for file in os.listdir(cache_folder):
    if file.endswith(".pt"):
        file_path = os.path.join(cache_folder, file)
        graph = torch.load(file_path, weights_only=False) # ⚠️ Explicitly disable weights_only
        graph_data.append(graph)

print(f"✅ Loaded {len(graph_data)} graph files from cache.")

# ✅ Fraction of the dataset to use
fraction = 0.75 # Use 100% of the dataset to reduce RAM usage
subset_size = int(len(graph_data) * fraction)
graph_data = graph_data[:subset_size] # Use only the fraction

# ✅ Split dataset into training and testing (80/20 split)
train_size = int(0.8 * len(graph_data))
test_size = len(graph_data) - train_size
train_data = graph_data[:train_size]
test_data = graph_data[train_size:]

# ✅ Create DataLoader
batch_size = 8 # Use smaller batch size to prevent RAM exhaustion
train_loader = DataLoader(train_data, batch_size=batch_size, shuffle=True)
test_loader = DataLoader(test_data, batch_size=batch_size, shuffle=False)

# ✅ Confirm batch shapes
for batch in train_loader:
    print(f"🔥 Batch Node Shape: {batch.x.shape}")
    print(f"🔥 Batch Edge Shape: {batch.edge_index.shape}")
    print(f"🔥 Labels Shape: {batch.y.shape}")
    break

🔄 ✅ Loaded 50 graph files from cache.
/usr/local/lib/python3.10/dist-packages/torch_geometric/deprecation.py:26: UserWarning: 'data.DataLoader' is deprecated, use 'loader
warnings.warn(out)
🔥 Batch Node Shape: torch.Size([11015891, 5])
🔥 Batch Edge Shape: torch.Size([2, 110158910])
🔥 Labels Shape: torch.Size([8])

```

```

# for batch in train_loader:
#     print("🔥 Batch Features Shape:", batch.x.shape) # Node features
#     print("🔥 Batch Edges Shape:", batch.edge_index.shape) # Edges
#     print("🔥 Labels:", batch.y) # Target labels
#     break

import collections

# ✅ Check label distribution in the entire dataset
all_labels = []
for batch in train_loader:
    all_labels.extend(batch.y.tolist())

print("🔥 Label Distribution:", collections.Counter(all_labels))

🔄 🔥 Label Distribution: Counter({1: 16, 0: 13})

import os
import time
import gc
import torch
import torch_geometric
from torch_geometric.nn import global_mean_pool # Import required modules

```

```

from torch_geometric.data import DataLoader

# -----
# ✅ Model Initialization
# -----
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
model = DGCNN(num_features=5).to(device)
optimizer = torch.optim.Adam(model.parameters(), lr=0.001)
criterion = torch.nn.CrossEntropyLoss()

# -----
# ✅ Training Configuration
# -----
num_epochs = 35
best_acc = 0.0
early_stopping_patience = 5
epochs_no_improve = 0
model_dir = "/kaggle/working/"
os.makedirs(model_dir, exist_ok=True)

# ✅ Memory-Efficient DataLoader
batch_size = 1 # Micro-batching to reduce VRAM consumption
num_accumulation_steps = 8 # Gradient accumulation over 8 micro-batches

train_loader = DataLoader(train_data, batch_size=batch_size, shuffle=True, pin_memory=True, num_workers=4)
test_loader = DataLoader(test_data, batch_size=batch_size, shuffle=False, pin_memory=True, num_workers=4)

# ✅ Metrics Tracking
train_losses, test_losses, train_accuracies, test_accuracies, epoch_times = [], [], [], [], []

# ✅ Mixed Precision + Gradient Accumulation
scaler = torch.cuda.amp.GradScaler()

# -----
# ✅ Training & Evaluation Loop
# -----
print("\n🔥 Training Started\n")

for epoch in range(1, num_epochs + 1):
    start_time = time.time()
    model.train()
    total_loss = 0.0

    # ✅ Gradient Accumulation Variables
    optimizer.zero_grad(set_to_none=True)

    # ✅ Micro-Batching & Accumulation
    for i, batch in enumerate(train_loader):
        batch = batch.to(device)

        with torch.cuda.amp.autocast():
            output = model(batch)
            loss = criterion(output, batch.y)

        # ✅ Accumulate Gradients
        scaler.scale(loss / num_accumulation_steps).backward()

        # ✅ Optimizer Step After Accumulating Gradients
        if (i + 1) % num_accumulation_steps == 0 or i == len(train_loader) - 1:
            torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
            scaler.step(optimizer)
            scaler.update()
            optimizer.zero_grad(set_to_none=True)

    total_loss += loss.item()

    # ✅ Memory Cleanup
    del batch
    torch.cuda.empty_cache()
    gc.collect()

# -----
# ✅ Evaluation
# -----
model.eval()
train_correct = 0
test_correct = 0

with torch.no_grad():
    for batch in train_loader:
        batch = batch.to(device)
        pred = model(batch).argmax(dim=1)

```

```

train_correct += (pred == batch.y).sum().item()

for batch in test_loader:
    batch = batch.to(device)
    pred = model(batch).argmax(dim=1)
    test_correct += (pred == batch.y).sum().item()

train_acc = train_correct / len(train_data)
test_acc = test_correct / len(test_data)

epoch_time = time.time() - start_time
epoch_times.append(epoch_time)

print(f"🔥 Epoch {epoch}/{num_epochs} | Loss: {total_loss:.4f} | Train Acc: {train_acc:.4f} | Test Acc: {test_acc:.4f} | Time: {epoch_time:.2f}s")

# ✅ Save Model After Each Epoch
torch.save(model.state_dict(), os.path.join(model_dir, f"model_epoch_{epoch}.pt"))

# ✅ Early Stopping
if test_acc > best_acc:
    best_acc = test_acc
    epochs_no_improve = 0
    print(f"🔥 Best model updated at epoch {epoch} with Test Acc: {best_acc:.4f}")
    with open(os.path.join(model_dir, f"model_epoch_{epoch}.pt"), "w") as f:
        f.write(f"Model saved at epoch {epoch} with Test Acc: {test_acc:.4f}")
else:
    epochs_no_improve += 1
    if epochs_no_improve >= early_stopping_patience:
        print(f"🔥 Early stopping triggered!")
        break

print("\n✅ Training Completed:")
print(f"🔥 Best Accuracy: {best_acc:.4f}")
print(f"🔥 Model saved at epoch {final_epoch}")

```



🔥 **Training Started:**

🔥 Epoch 1/35	Loss: 1.3329	Train Acc: 0.5715	Test Acc: 0.5283	Time: 380.53s
🔥 Epoch 2/35	Loss: 1.3006	Train Acc: 0.6032	Test Acc: 0.5682	Time: 347.95s
🔥 Epoch 3/35	Loss: 1.2504	Train Acc: 0.6350	Test Acc: 0.5947	Time: 308.28s
🔥 Epoch 4/35	Loss: 1.1717	Train Acc: 0.6616	Test Acc: 0.6238	Time: 320.62s
🔥 Epoch 5/35	Loss: 1.0517	Train Acc: 0.6789	Test Acc: 0.6602	Time: 359.80s
🔥 Epoch 6/35	Loss: 1.0101	Train Acc: 0.6985	Test Acc: 0.6939	Time: 310.85s
🔥 Epoch 7/35	Loss: 0.9357	Train Acc: 0.7227	Test Acc: 0.7277	Time: 369.44s
🔥 Epoch 8/35	Loss: 0.8668	Train Acc: 0.7488	Test Acc: 0.7676	Time: 335.45s
🔥 Epoch 9/35	Loss: 0.7580	Train Acc: 0.7725	Test Acc: 0.7525	Time: 377.42s
🔥 Epoch 10/35	Loss: 0.7224	Train Acc: 0.7954	Test Acc: 0.7468	Time: 365.50s
🔥 Epoch 11/35	Loss: 0.6543	Train Acc: 0.8089	Test Acc: 0.7480	Time: 370.31s
🔥 Epoch 12/35	Loss: 0.6318	Train Acc: 0.8080	Test Acc: 0.7556	Time: 359.11s
🔥 Best model updated at epoch 12 with Test Acc: 0.7556				
🔥 Epoch 13/35	Loss: 0.5569	Train Acc: 0.8075	Test Acc: 0.7636	Time: 352.38s
🔥 Best model updated at epoch 13 with Test Acc: 0.7636				
🔥 Epoch 14/35	Loss: 0.5569	Train Acc: 0.8086	Test Acc: 0.7530	Time: 351.79s
🔥 Epoch 15/35	Loss: 0.5829	Train Acc: 0.8007	Test Acc: 0.7641	Time: 321.30s
🔥 Best model updated at epoch 15 with Test Acc: 0.7641				
🔥 Epoch 16/35	Loss: 0.5569	Train Acc: 0.8089	Test Acc: 0.7644	Time: 308.21s
🔥 Best model updated at epoch 16 with Test Acc: 0.7644				
🔥 Epoch 17/35	Loss: 0.5569	Train Acc: 0.7996	Test Acc: 0.7568	Time: 388.62s
🔥 Epoch 18/35	Loss: 0.5812	Train Acc: 0.8125	Test Acc: 0.7568	Time: 387.88s
🔥 Epoch 19/35	Loss: 0.5569	Train Acc: 0.8107	Test Acc: 0.7577	Time: 300.92s
🔥 Epoch 20/35	Loss: 0.5569	Train Acc: 0.8108	Test Acc: 0.7570	Time: 354.40s
🔥 Epoch 21/35	Loss: 0.5680	Train Acc: 0.8095	Test Acc: 0.7584	Time: 333.36s
🔥 Epoch 22/35	Loss: 0.5569	Train Acc: 0.8167	Test Acc: 0.7579	Time: 324.09s
🔥 Epoch 23/35	Loss: 0.5569	Train Acc: 0.7992	Test Acc: 0.7693	Time: 353.66s
🔥 Best model updated at epoch 23 with Test Acc: 0.7693				
🔥 Epoch 24/35	Loss: 0.5708	Train Acc: 0.7959	Test Acc: 0.7541	Time: 334.85s
🔥 Epoch 25/35	Loss: 0.5569	Train Acc: 0.8108	Test Acc: 0.7553	Time: 373.42s
🔥 Epoch 26/35	Loss: 0.5569	Train Acc: 0.8064	Test Acc: 0.7487	Time: 390.37s
🔥 Epoch 27/35	Loss: 0.5754	Train Acc: 0.7995	Test Acc: 0.7585	Time: 366.57s
🔥 Epoch 28/35	Loss: 0.5569	Train Acc: 0.7970	Test Acc: 0.7601	Time: 398.39s
🔥 Early stopping triggered!				

✅ Training Completed:
 🔄 Best Accuracy: 0.7693
 🔥 Model saved at epoch 28

```
import matplotlib.pyplot as plt
```

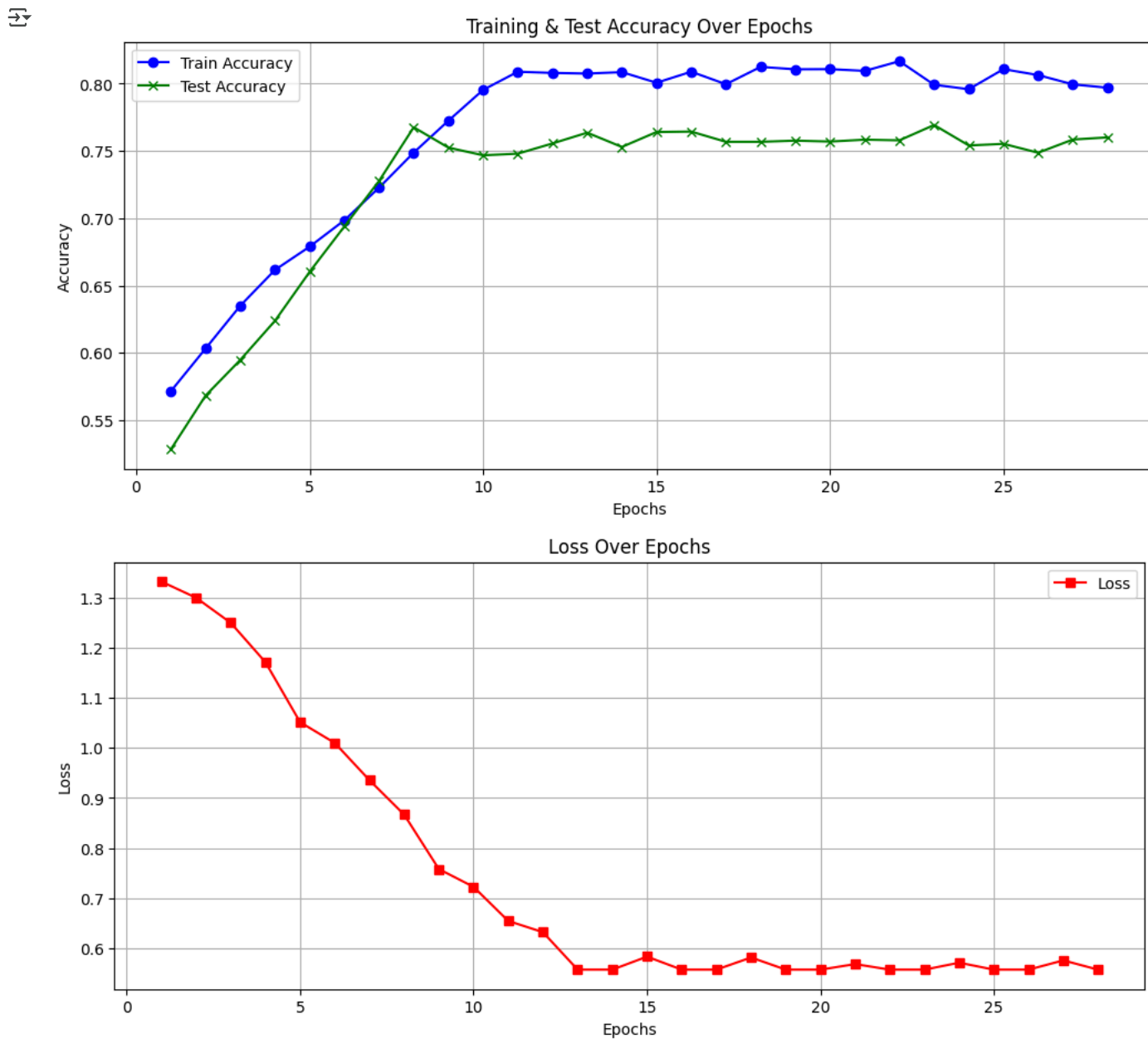
```

# ✅ Plotting Training and Test Accuracy
plt.figure(figsize=(12, 5))
plt.plot(epochs, train_acc, label='Train Accuracy', marker='o', color='blue')
plt.plot(epochs, test_acc, label='Test Accuracy', marker='x', color='green')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')

```

```
plt.title('Training & Test Accuracy Over Epochs')
plt.legend()
plt.grid(True)
plt.show()

# ☒ Plotting Loss
plt.figure(figsize=(12, 5))
plt.plot(epochs, loss, label='Loss', marker='s', color='red')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.title('Loss Over Epochs')
plt.legend()
plt.grid(True)
plt.show()
```



```
# -----
# ☒ Import Libraries
# -----
import os
import torch
import numpy as np
import pandas as pd
from torch_geometric.data import Data
from scipy.spatial import cKDTree
from sklearn.preprocessing import MinMaxScaler
import h5py
from tqdm import tqdm

# -----
# ☒ Constants
# -----
```

```

HDF5_FILE_PATH = '/kaggle/input/quark-gluon-ds/quark-gluon_data-set_n139306.hdf5'
MODEL_PATH = '/kaggle/input/clf-model/best_model_clf.pt'
K = 5 # Nearest neighbors for kNN graph
NUM_IMAGES = 100 # Last 100 images for inference

# -----
# ✅ Load the Trained Model
# -----
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

# ✅ Model Definition (DGCNN architecture)
class DGCNN(torch.nn.Module):
    def __init__(self, num_features):
        super(DGCNN, self).__init__()
        self.conv1 = torch.nn.Linear(num_features, 128)
        self.conv2 = torch.nn.Linear(128, 64)
        self.fc = torch.nn.Linear(64, 2) # Binary classification: Quark vs Gluon

    def forward(self, data):
        x, edge_index = data.x, data.edge_index
        x = torch.relu(self.conv1(x))
        x = torch.relu(self.conv2(x))
        x = torch.mean(x, dim=0) # Global pooling
        x = self.fc(x)
        return torch.softmax(x, dim=0)

# ✅ Initialize and load model
model = DGCNN(num_features=5).to(device)
model.load_state_dict(torch.load(MODEL_PATH, map_location=device))
model.eval()
print("✅ Model loaded successfully!")

# -----
# ✅ Preprocessing Functions
# -----
def knn_graph(x, k=5):
    """
    Constructs a kNN graph from node features.
    Args:
        x (torch.Tensor): Node features [num_nodes, num_features].
        k (int): Number of neighbors.
    Returns:
        edge_index (torch.Tensor): [2, num_edges].
    """
    x_np = x.cpu().numpy()
    tree = cKDTree(x_np)

    _, neighbors = tree.query(x_np, k=k + 1) # k+1 to include self
    edge_index = []

    for i, n in enumerate(neighbors):
        for j in n[1:]: # Skip self-loop
            edge_index.append([i, j])
            edge_index.append([j, i])

    edge_index = torch.tensor(edge_index, dtype=torch.long).T
    return edge_index

def preprocess_image(img):
    """
    Preprocesses an image into point cloud and constructs kNN graph.
    Args:
        img (np.array): The 3D image array.
    Returns:
        Data: PyG Data object with the graph.
    """
    # ✅ Extract non-zero pixels
    non_zero_mask = np.any(img != 0, axis=-1)
    y_coords, x_coords = np.where(non_zero_mask)

    # ✅ Extract intensity features (ECAL, HCAL, Tracks)
    ecal = img[non_zero_mask][:, 0]
    hcal = img[non_zero_mask][:, 1]
    tracks = img[non_zero_mask][:, 2]

    # ✅ Normalize features
    features = np.column_stack((ecal, hcal, tracks))
    scaler = MinMaxScaler()
    features_normalized = scaler.fit_transform(features)

    # ✅ Create Point Cloud DataFrame
    point_cloud = pd.DataFrame({

```