

The dataset has four main elements:

X\_jets – Three-channel images, all 125 × 125 pixels in size.

m0 – Stands for mass.

pt – Stands for transverse momentum.

y – Refers to labels.

State-of-the-art Jet Image Taggers make use of a 4-tuple representation when performing classification.

For this particular screening task, I am implementing a Variational Autoencoder (VAE) to examine and process Jet Images as needed.

```
In [1]: import h5py
import numpy as np
import torch
import torch.nn as nn
import torch.optim as optim
import matplotlib.pyplot as plt
from torch.utils.data import Dataset, DataLoader
from torch.cuda.amp import autocast, GradScaler
from sklearn.metrics import confusion_matrix
from torch.distributions import Normal
from tqdm import tqdm
from sklearn.metrics_selection import train_test_split

In [2]: device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
print(f"Using device: {device}")

Using device: cuda

In [3]: file_path = "haggle/input/autocoder-data/quark-gluon-data-set_v139306.hdf5"
def explore_hdf5(file):
    with h5py.File(file, 'r') as f:
        print("Dataset keys:", list(f.keys()))
        for key in f.keys():
            print(f"Shape of {key}: {f[key].shape}")

explore_hdf5(file_path)

Keys in dataset: ['X_jets', 'm0', 'pt', 'y']
Shape of X_jets: (139306, 125, 125, 3)
Shape of m0: (139306,)
Shape of pt: (139306,)
Shape of y: (139306,)

In [4]: def load_data(file_name, sample_size):
    with h5py.File(file_name, 'r') as f:
        print("Dataset keys:", list(f.keys()))
        print("Total images:", len(f['X_jets']))
        print("Using sample size:", f['X_jets'].shape[0])
        return np.array(f['X_jets'][:sample_size]), np.array(f['y'][:sample_size])

X, y = load_data(file_name, 10000)

Dataset keys: ['X_jets', 'm0', 'pt', 'y']
Total images: 139306
Using dimensions: (125, 125, 3)
```

This visualize\_channels function will plot and compare a multi-channel jet image by showing:

The combined image (all three channels combined) on top. Single channel images (Track, ECAL, and HCAL) underneath it, with 'viridis' colormap for better visualization.

```
In [5]: def visualize_channels(image_data, sample_index=129):
    X_sample = image_data[sample_index]
    fig, axes = plt.subplots(4, 1, figsize=(10, 20))
    # Standardizing the distribution across all channels
    axes[0].set_title('Combined')
    for i in range(3):
        img = axes[i+1].imshow(X_sample[:, i, :, :], cmap='viridis', vmin=-0.5, vmax=2.0, interpolation='nearest')
        axes[i+1].set_title(channel_labels[i])
    fig.colorbar(img, ax=axes[i+1], shrink=0.5)
    plt.show()

visualize_channels(X)
```

This function intensity calculates and prints statistical information regarding the intensity values over various channels in the dataset.

Traverses all three channels (Track, ECAL, HCAL) to calculate:

Maximum and Minimum intensity values Mean intensity value Standard deviation of intensity distribution Calculates an overall intensity image by adding up the three channels and studying its statistical properties.

```
In [6]: def intensity(X, y):
    for i, channel_name in enumerate(['Track', 'ECAL', 'HCAL']):
        print(f"Max value of intensity along {channel_name} channel: {np.max(X[:, i, :, :])}, Min value: {np.min(X[:, i, :, :])}")
        print(f"Mean intensity value along {channel_name} channel: {np.mean(X[:, i, :, :])}")
        print(f"Standard Deviation: {np.std(X[:, i, :, :])}")

    combined_dataset = np.sum(X, axis=-1, keepdims=True)
    print("Max value of intensity in combined channel image:", np.max(combined_dataset[:, :, :]), "Min value:", np.min(combined_dataset[:, :, :]))
    print("Mean intensity value in combined channel image:", np.mean(combined_dataset[:, :, :]))
    print("Standard Deviation:", np.std(combined_dataset[:, :, :]), "\n")

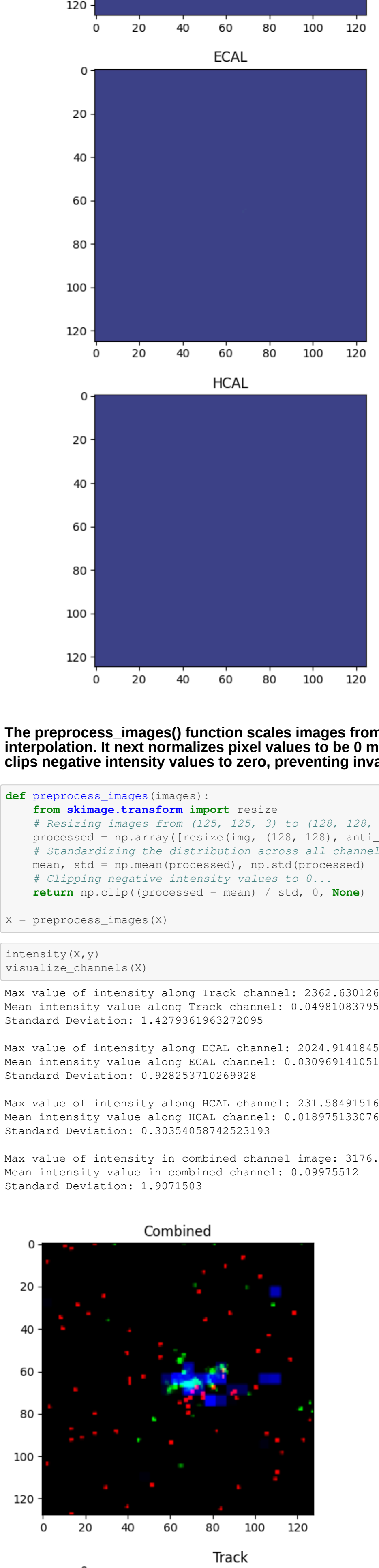
intensity(X, y)

Max value of intensity along Track channel: 2.0881052017211891, Min value: 0.0
Mean intensity value along Track channel: 0.841042679501697e-05
Standard Deviation: 0.008757636211812496

Max value of intensity along ECAL channel: 0.35979661345481857, Min value: 0.0
Mean intensity value along ECAL channel: 4.968207210962086e-05
Standard Deviation: 0.002107673790305853

Max value of intensity along HCAL channel: 0.35979661345481857, Min value: 0.0
Mean intensity value along HCAL channel: 3.21495158989644e-05
Standard Deviation: 0.0005132272497914732

Max value of intensity in combined channel image: 12.1562605 Min value: 0.0
Mean intensity value in combined channel: 0.0015928746
Standard Deviation: 0.004358586
```



The preprocess\_images() function scales images from (125, 125, 3) to (128, 128, 3) for uniformity with bilinear interpolation. It next normalizes pixel values to be 0 mean and 1 standard deviation, stabilizing the training. It lastly clips negative intensity values to zero, preventing invalid pixel values and lowering noise.

```
In [8]: def preprocess_images(images):
    # Resizing images from (125, 125, 3) to (128, 128, 3)
    processed = np.array([resize(img, (128, 128), anti_aliasing=True) for img in images], dtype=np.float32)
    # Standardizing the distribution across all channels
    mean, std = np.mean(processed), np.std(processed)
    # Clipping negative intensity values to 0
    return np.clip(processed - mean) / std, 0, None

X = preprocess_images(X)

In [9]: intensity(X, y)
visualize_channels(X)

Max value of intensity along Track channel: 2.362163012693125, Min value: 0.0
Mean intensity value along Track channel: 0.0491808759428276
Standard Deviation: 1.4273181561270395

Max value of intensity along ECAL channel: 2.0241118570311325, Min value: 0.0
Mean intensity value along ECAL channel: 0.03098914105117321
Standard Deviation: 0.028233710249928

Max value of intensity along HCAL channel: 2.3115819151611028, Min value: 0.0
Mean intensity value along HCAL channel: 0.02897131079610338
Standard Deviation: 0.30354058742523193

Max value of intensity in combined channel image: 13.761876 Min value: 0.0
Mean intensity value in combined channel: 0.03975212
Standard Deviation: 1.9071503
```



```
In [10]: class Reparameterization(nn.Module):
    def __init__(self, mean, log_variance):
        batch_size, latent_dim = mean.shape
        epsilon = Normal(0, 1).sample((batch_size, latent_dim)).to(mean.device)
        return mean + torch.exp(0.5 * log_variance) * epsilon

In [11]: class VAE_autoencoder(nn.Module):
    def __init__(self, embedding_dim):
        super().__init__()
        self.encoder = nn.Sequential(
            nn.Conv2d(3, 32, kernel_size=3, stride=2, padding=1),
            nn.BatchNorm2d(32),
            nn.LeakyReLU(0.2),
            nn.Conv2d(32, 64, kernel_size=3, stride=2, padding=1),
            nn.BatchNorm2d(64),
            nn.LeakyReLU(0.2),
            nn.Conv2d(64, 128, kernel_size=3, stride=2, padding=1),
            nn.BatchNorm2d(128),
            nn.LeakyReLU(0.2),
            nn.Conv2d(128, 256, kernel_size=3, stride=2, padding=1),
            nn.BatchNorm2d(256),
            nn.LeakyReLU(0.2),
        )
        self.mean_layer = nn.Linear(6384, embedding_dim)
        self.logvar = nn.Linear(6384, embedding_dim)
        self.reparam = Reparameterization()

        self.decoder = nn.Sequential(
            nn.Linear(embedding_dim, 6384),
            nn.ConvTranspose2d(256, 128, kernel_size=3, stride=2, padding=1, output_padding=1),
            nn.BatchNorm2d(128),
            nn.LeakyReLU(0.2),
            nn.ConvTranspose2d(128, 64, kernel_size=3, stride=2, padding=1, output_padding=1),
            nn.BatchNorm2d(64),
            nn.LeakyReLU(0.2),
            nn.ConvTranspose2d(64, 32, kernel_size=3, stride=2, padding=1, output_padding=1),
            nn.BatchNorm2d(32),
            nn.LeakyReLU(0.2),
            nn.ConvTranspose2d(32, 3, kernel_size=3, stride=2, padding=1, output_padding=1),
            nn.Tanh()
        )

    def forward(self, x, embedding_dim):
        z = self.encoder(x)
        x = x.reshape(x.size()[0], -1)
        mu = self.mean_layer(x)
        logvar = self.logvar(x)
        z = self.reparam(mu, logvar)
        z = self.decoder(z)
        z = z.view(z.size()[0], 256, 8, 8)

        # Continue with the decoder
        x_rec = self.decoder(z)
        x_rec = (x_rec + 1) / 2
        return x_rec, mu, logvar

In [12]: class VAEloss(nn.Module):
    def __init__(self):
        super().__init__()
        self.kl_loss = nn.BCELoss()

    def forward(self, x_rec, mu, logvar):
        recon_loss = self.bce_loss(x_rec, x)
        kl_loss = -0.5 * torch.sum(1 - logvar - mu.pow(2) - logvar.exp(), dim=-1).mean()
        return 500 * recon_loss + kl_loss

In [13]: trackMax = np.max(X[:, 1, :, :])
ecalMax = np.max(X[:, 2, :, :])
hcalMax = np.max(X[:, 3, :, :])

X[:, :, 1, 0] = X[:, :, 1, 0] / trackMax
X[:, :, 1, 1] = X[:, :, 1, 1] / ecalMax
X[:, :, 1, 2] = X[:, :, 1, 2] / hcalMax

In [14]: intensity(X, y)

Max value of intensity along Track channel: 1.0, Min value: 0.0
Mean intensity value along Track channel: 2.10827911359956e-05
Standard Deviation: 0.00604364225184731

Max value of intensity along ECAL channel: 1.0, Min value: 0.0
Mean intensity value along ECAL channel: 8.193597022909671e-05
Standard Deviation: 0.0011070791516521729

Max value of intensity in combined channel image: 1.6550478 Min value: 0.0
Mean intensity value in combined channel: 0.001181283
Standard Deviation: 0.016820203
```

Loading these datasets using standard PyTorch dataloader ->

```
In [15]: X_jets_train, X_jets_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
X_jets_train, X_jets_val, y_train, y_val = train_test_split(X_jets_train, y_train, test_size=0.25, random_state=42)

In [16]: class JetDataset(Dataset):
    def __init__(self, X, y):
        self.X = X
        self.y = y

    def __getitem__(self, index):
        image = self.X[index].permute(2, 0, 1)
        return image, self.y[index]

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

train_dataset = JetDataset(torch.from_numpy(X_jets_train), torch.from_numpy(y_train))
val_dataset = JetDataset(torch.from_numpy(X_jets_val), torch.from_numpy(y_val))
test_dataset = JetDataset(torch.from_numpy(X_jets_test), torch.from_numpy(y_test))

In [17]: from torch import nn

In [18]: device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
model = VAE_autoencoder(embedding_dim=2048).to(device)
optimizer = optim.Adam(model.parameters())
criterion = VAELoss()

# Data loaders
train_loader = DataLoader(train_dataset, batch_size=128, shuffle=False)
val_loader = DataLoader(val_dataset, batch_size=128, shuffle=False)

# Training loop
num_epochs = 20
train_losses = []
val_losses = []
latent_representations = []
labels = []

for epoch in range(num_epochs):
    model.train()
    train_loss = 0
    for batch_idx, (data, _) in enumerate(val_loader):
        for batch_idx, (data, _) in enumerate(train_loader, desc=f"Epoch {epoch}"):
            data = data.to(device)
            optimizer.zero_grad()
            recon_batch, mu, logvar = model(data, embedding_dim=2048)
            loss = criterion(recon_batch, data, mu, logvar)
            loss.backward()
            train_loss += loss.item()
            train_loss /= len(train_loader)
            optimizer.step()
            train_loss.append(train_loss / len(train_loader))

    model.eval()
    val_loss = 0
    with torch.no_grad():
        for batch_idx, (data, _) in enumerate(val_loader):
            data = data.to(device)
            recon_batch, mu, logvar = model(data, embedding_dim=2048)
            loss = criterion(recon_batch, data, mu, logvar)
            val_loss += loss.item()
            val_loss /= len(val_loader)

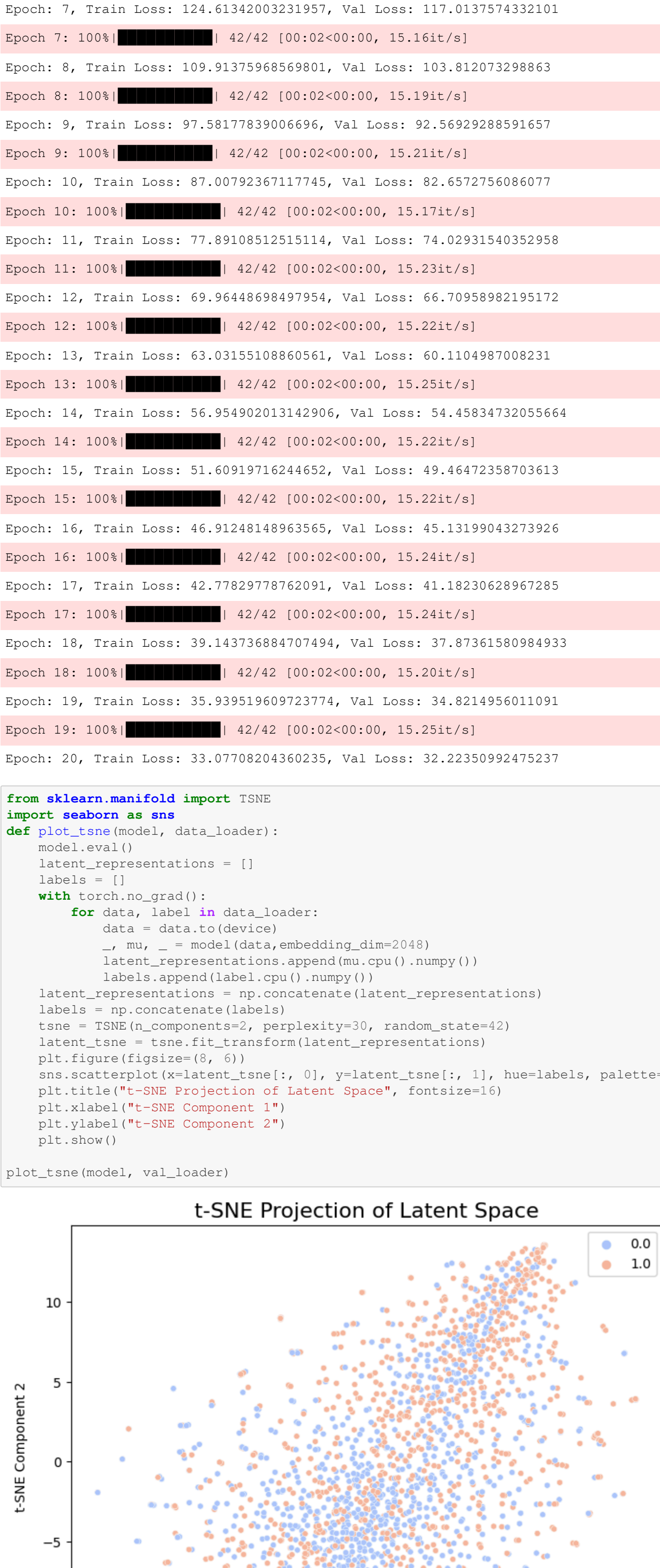
    print(f"Epoch: {epoch+1}, Train Loss: {train_loss:.4f}, Val Loss: {val_loss:.4f}")

Epoch 0: 100% ██████████ 42/42 [00:03:00:00, 11.001e/s]
Epoch: 1, Train Loss: 435.89149402436755, Val Loss: 337.7340300873326
Epoch: 2, Train Loss: 311.492023732824, Val Loss: 248.3485761369978
Epoch: 3, Train Loss: 244.51296960739862, Val Loss: 216.306817141202
Epoch: 4, Train Loss: 197.7042937505813, Val Loss: 178.73394557407923
Epoch: 5, Train Loss: 165.86970919964482, Val Loss: 152.979127066476
Epoch: 6, Train Loss: 142.6474903784967, Val Loss: 133.0473597865515
Epoch: 7, Train Loss: 124.61342003231957, Val Loss: 117.0137573323101
Epoch: 8, Train Loss: 109.91375988569801, Val Loss: 103.8160732398863
Epoch: 9, Train Loss: 97.58178390064696, Val Loss: 92.5692988591657
Epoch: 10, Train Loss: 87.002367117745, Val Loss: 82.657256806077
Epoch: 11, Train Loss: 77.89108512515114, Val Loss: 74.02931540352958
Epoch: 12, Train Loss: 69.9644869497954, Val Loss: 66.70989882195172
Epoch: 13, Train Loss: 63.0315108860561, Val Loss: 60.110498706231
Epoch: 14, Train Loss: 56.954902013142906, Val Loss: 54.4583472055664
Epoch: 15, Train Loss: 51.60919716244652, Val Loss: 49.464723587031613
Epoch: 16, Train Loss: 46.91248148963565, Val Loss: 45.13199043273926
Epoch: 17, Train Loss: 42.7782978782091, Val Loss: 41.18230628967285
Epoch: 18, Train Loss: 39.143736884707494, Val Loss: 37.873614508984933
Epoch: 19, Train Loss: 35.93951960923774, Val Loss: 34.821456011091
Epoch: 20, Train Loss: 33.07708204360235, Val Loss: 32.22350992475237
```

```
In [19]: from sklearn.manifold import TSNE
import seaborn as sns

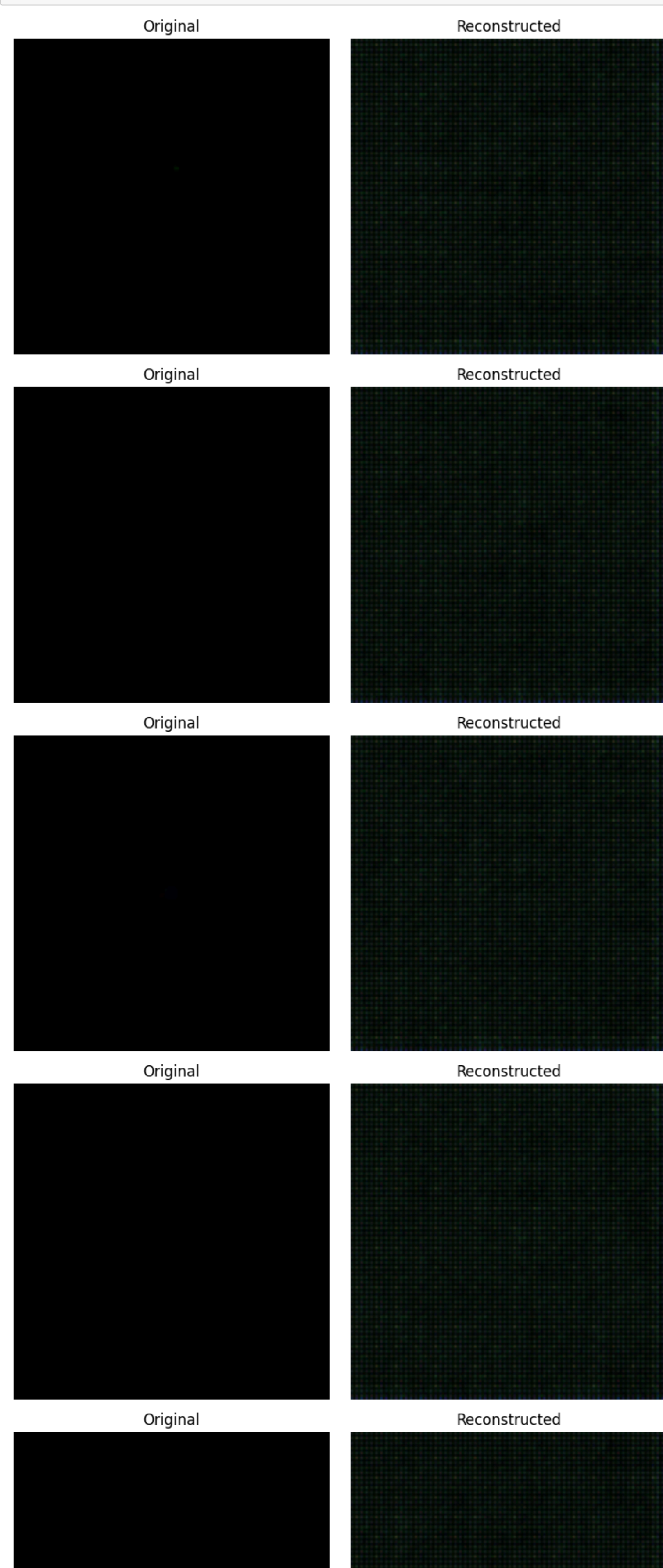
def plot_latent_space(model, data_loader):
    model.eval()
    latent_representations = []
    labels = []
    with torch.no_grad():
        for data, label in data_loader:
            data = data.to(device)
            mu, _ = model(data, embedding_dim=2048)
            latent_representations.append(mu.cpu().numpy())
            labels.append(label.cpu().numpy())
    latent_representations = np.concatenate(latent_representations)
    labels = np.concatenate(labels)
    t_sne = TSNE(n_components=2, perplexity=30, random_state=42)
    latent_t_sne = t_sne.fit_transform(latent_representations)
    plt.figure(figsize=(8, 8))
    sns.scatterplot(latent_t_sne[:, 0], latent_t_sne[:, 1], hue=labels, palette="coolwarm", s=20)
    plt.title(f"t-SNE Projection of Latent Space", fontsize=16)
    plt.xlabel(f"t-SNE Component 1")
    plt.ylabel(f"t-SNE Component 2")
    plt.show()

plot_latent_space(model, val_loader, num_samples=5)
```

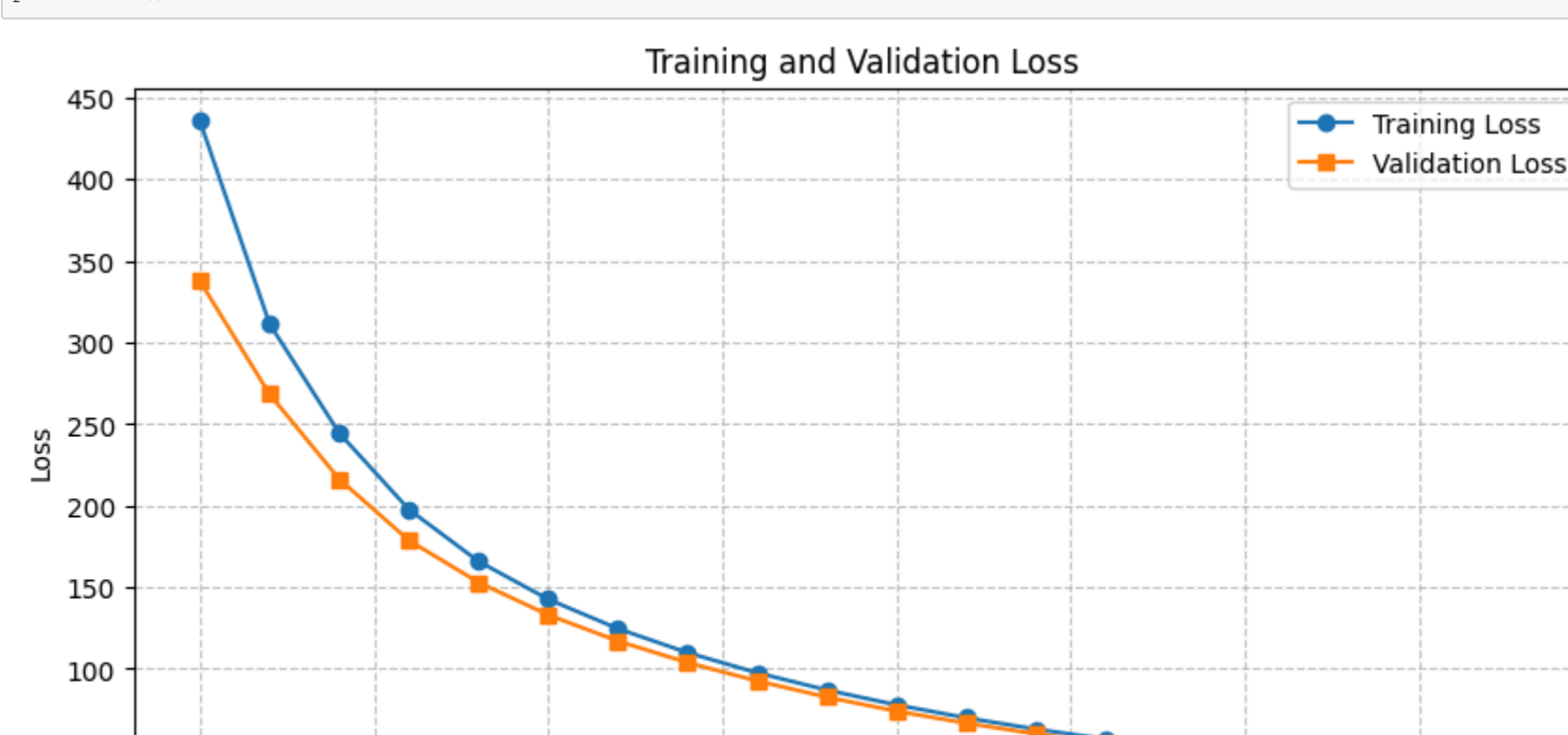


```
In [20]: def visualize_reconstruction(model, data_loader, num_samples=5):
    model.eval()
    recon_batch, _ = next(iter(data_loader))
    recon_batch = recon_batch.to(device)
    with torch.no_grad():
        _, mu = model(recon_batch, embedding_dim=2048)
        reconstructed = torch.nn.functional.sigmoid(mu)
        reconstructed = reconstructed.cpu().numpy().transpose(0, 2, 3, 1)
    fig, axes = plt.subplots(num_samples, 2, figsize=(8, num_samples * 4))
    for i in range(num_samples):
        axes[i, 0].imshow(recon_batch[i])
        axes[i, 1].imshow(recon_batch[i])
        axes[i, 0].set_title('Original')
        axes[i, 1].set_title('Reconstructed')
        axes[i, 0].axis('off')
        axes[i, 1].axis('off')
    plt.tight_layout()
    plt.show()

In [21]: visualize_comparison(model, val_loader, num_samples=5)
```



```
In [22]: plt.figure(figsize=(10, 5))
plt.plot(train_losses, marker='o', linestyle='-', label='Training Loss')
plt.plot(val_losses, marker='x', linestyle='-', label='Validation Loss')
plt.title('Training and Validation Loss')
plt.xlabel('Epochs')
plt.ylabel('Loss')
plt.legend()
plt.grid(True, linestyle='--', alpha=0.7)
plt.show()
```



```
In [ ]:
```