

Group 130

Members	Mail ID	Contribution %
Raghavendar V	2024aa05458@wilp.bits-pilani.ac.in	100%
Vikas Bansal	2024aa05457@wilp.bits-pilani.ac.in	100%
Vignesh P	2024aa05605@wilp.bits-pilani.ac.in	100%
Rajat Kumar Singh	2024aa05460@wilp.bits-pilani.ac.in	100%
Abhay Suresh Gunjarkar	2023ac05062@wilp.bits-pilani.ac.in	100%

Advanced Deep Learning - Assignment - 2

In this assignment, you will implement and compare four generative modeling approaches for face modification and synthesis:

1. Variational Autoencoder (VAE)
2. β -Variational Autoencoder (β -VAE)
3. Vector Quantized VAE (VQ-VAE) with PixelCNN prior for latent sampling
4. Generative Adversarial Network (GAN)

The goal is to understand how different latent representations influence reconstruction quality, attribute manipulation, and sample realism in face images.

Use provisioned WILP lab infrastructure! Dataset Use CelebA dataset.

Preprocessing • Resize images to 64×64 or 128×128 • Normalize pixel values to [0,1] • Optionally select attributes (e.g., Smiling, Male, Eyeglasses)

Task

Downloading and preprocessing the CelebA dataset, resizing images to 64x64 and normalizing pixel values to [0, 1] as specified in the assignment.

Prepare CelebA Dataset

Subtask:

Download and preprocess the CelebA dataset (resize to 64x64, normalize to [0, 1]).

Importing the necessary libraries, define the configuration, create the transformation pipeline, load the CelebA dataset using torchvision, create a DataLoader, and then visualize a batch of images while printing their statistics to verify the preprocessing.

```

import torch
import torchvision
import torchvision.transforms as transforms
import torch.utils.data as data
import matplotlib.pyplot as plt
import numpy as np

# 2. Define configuration variables
image_size = 64
batch_size = 64

# 3. Create transformation pipeline
transform = transforms.Compose([
    transforms.Resize((image_size, image_size)),
    transforms.ToTensor(),
])

# 4. Initialize CelebA dataset
# Note: This attempts to download the dataset. If the Google Drive
# quota is exceeded,
# you might need to manually download or use a different source.
try:
    train_dataset = torchvision.datasets.CelebA(root='./data',
        split='train', download=True, transform=transform)
except RuntimeError as e:
    print(f"Download failed: {e}. Checking if data already exists...")
    train_dataset = torchvision.datasets.CelebA(root='./data',
        split='train', download=False, transform=transform)

# 5. Create DataLoader
dataloader = data.DataLoader(train_dataset, batch_size=batch_size,
    shuffle=True)

# 6. Retrieve a single batch
real_batch = next(iter(dataloader))
images = real_batch[0]

# 8. Print shape and value ranges
print(f"Batch shape: {images.shape}")
print(f"Pixel values range: [{images.min():.4f}, {images.max():.4f}]")

# 7. Visualize a subset of images
plt.figure(figsize=(8, 8))
plt.axis("off")
plt.title("Training Images")
# Create grid
grid_img = torchvision.utils.make_grid(images[:16], padding=2,
    normalize=False)
# Transpose for plotting (C, H, W) -> (H, W, C)

```

```
plt.imshow(np.transpose(grid_img.numpy(), (1, 2, 0)))
plt.show()

Downloading...
From (original): https://drive.google.com/uc?
id=0B7EVK8r0v71pZjFTYXZWM3FlRnM
From (redirected): https://drive.usercontent.google.com/download?
id=0B7EVK8r0v71pZjFTYXZWM3FlRnM&confirm=t&uuid=32a02ccc-e6fe-4f71-
bcb9-65fb0b0380f1
To: /content/data/celeba/img_align_celeba.zip
100%|██████████| 1.44G/1.44G [00:21<00:00, 67.8MB/s]
Downloading...
From: https://drive.google.com/uc?id=0B7EVK8r0v71pblRyaVFSWGxPY0U
To: /content/data/celeba/list_attr_celeba.txt
100%|██████████| 26.7M/26.7M [00:00<00:00, 67.1MB/s]
Downloading...
From: https://drive.google.com/uc?id=1_ee_0u7vcNL0fNLegJRHmolfH5ICW-XS
To: /content/data/celeba/identity_CelebA.txt
100%|██████████| 3.42M/3.42M [00:00<00:00, 177MB/s]
Downloading...
From: https://drive.google.com/uc?id=0B7EVK8r0v71pbThiMVRxWXZ4dU0
To: /content/data/celeba/list_bbox_celeba.txt
100%|██████████| 6.08M/6.08M [00:00<00:00, 37.0MB/s]
Downloading...
From: https://drive.google.com/uc?id=0B7EVK8r0v71pd0FJY3Blby1HUTQ
To: /content/data/celeba/list_landmarks_align_celeba.txt
100%|██████████| 12.2M/12.2M [00:00<00:00, 53.5MB/s]
Downloading...
From: https://drive.google.com/uc?id=0B7EVK8r0v71pY0NSMzRuSXJEVkk
To: /content/data/celeba/list_eval_partition.txt
100%|██████████| 2.84M/2.84M [00:00<00:00, 237MB/s]

Batch shape: torch.Size([64, 3, 64, 64])
Pixel values range: [0.0000, 1.0000]
```

Training Images



Reasoning: The previous execution completed successfully, including the loading, and visualization.

Task

Derive the Evidence Lower Bound (ELBO) loss for the Variational Autoencoder (VAE) with a mathematical explanation. Then, implement the VAE architecture in PyTorch, including the reparameterization trick. Create a training loop to train the model on the CelebA dataset (located at `./data`), and implement functions to visualize reconstructions, perform latent interpolation between two faces, and demonstrate attribute modification.

Part A: VAE Implementation and Analysis

Subtask:

Derive ELBO, implement VAE in PyTorch, train on CelebA, and visualize results (reconstruction, interpolation, attribute manipulation).

1. Mathematical Derivation of ELBO for VAE

The goal of the Variational Autoencoder (VAE) is to maximize the likelihood of the data $p(x)$. Since the marginal likelihood involves an intractable integral over the latent space z (i.e., $p(x) = \int p(x|z)p(z)dz$), we instead maximize a lower bound on the log-likelihood, known as the **Evidence Lower Bound (ELBO)**.

We introduce an approximate posterior distribution $q_\phi(z|x)$ (the encoder) to approximate the true posterior $p_\theta(z|x)$.

Starting from the log-likelihood $\log p(x)$:

$$\log p(x) = \log \int p(x, z) dz = \log E_{q_\phi(z|x)} \left[\frac{p_\theta(x|z)p(z)}{q_\phi(z|x)} \right]$$

By applying Jensen's Inequality (since \log is a concave function, $\log E[y] \geq E[\log y]$):

$$\log p(x) \geq E_{q_\phi(z|x)} \left[\log \frac{p_\theta(x|z)p(z)}{q_\phi(z|x)} \right]$$

This term on the right-hand side is the **ELBO**.

1. **Reconstruction Term** $E_{q_\phi(z|x)} [\log p_\theta(x|z)]$: Encourages the decoder to reconstruct the input x accurately from the latent variable z .
2. **KL Divergence** D_{KL} : Regularizes the encoder by forcing the approximate posterior $q_\phi(z|x)$ to be close to the prior $p(z)$ (usually a standard normal distribution $N(0, I)$).

The VAE is trained by minimizing the negative ELBO:

$$L_{VAE} = -E_{q_\phi(z|x)} [\log p_\theta(x|z)] + D_{KL}$$

Reasoning: Implementing the VAE architecture in PyTorch as defined in the instructions, including the Encoder, Reparameterization trick, and Decoder.

```
import torch
import torch.nn as nn
import torch.nn.functional as F

class VAE(nn.Module):
    def __init__(self, latent_dim=128):
        super(VAE, self).__init__()
        self.latent_dim = latent_dim

        # Encoder
        self.encoder = nn.Sequential(
            nn.Conv2d(3, 32, kernel_size=4, stride=2, padding=1), # 64
            nn.BatchNorm2d(32),
            nn.LeakyReLU(0.2, inplace=True),
            nn.Conv2d(32, 64, kernel_size=4, stride=2, padding=1), #
            nn.BatchNorm2d(64),
            nn.LeakyReLU(0.2, inplace=True),
            nn.Conv2d(64, 128, kernel_size=4, stride=2, padding=1), #
            nn.BatchNorm2d(128),
            nn.LeakyReLU(0.2, inplace=True),
            nn.Conv2d(128, 256, kernel_size=4, stride=2, padding=1), #
            nn.BatchNorm2d(256),
            nn.LeakyReLU(0.2, inplace=True)
        )

        self.fc_mu = nn.Linear(256 * 4 * 4, latent_dim)
        self.fc_logvar = nn.Linear(256 * 4 * 4, latent_dim)

        # Decoder
        self.decoder_input = nn.Linear(latent_dim, 256 * 4 * 4)

        self.decoder = nn.Sequential(
            nn.ConvTranspose2d(256, 128, kernel_size=4, stride=2,
padding=1), # 4 -> 8
            nn.BatchNorm2d(128),
            nn.ReLU(inplace=True),
            nn.ConvTranspose2d(128, 64, kernel_size=4, stride=2,
padding=1), # 8 -> 16
            nn.BatchNorm2d(64),
```

```

        nn.ReLU(inplace=True),
        nn.ConvTranspose2d(64, 32, kernel_size=4, stride=2,
padding=1), # 16 -> 32
        nn.BatchNorm2d(32),
        nn.ReLU(inplace=True),
        nn.ConvTranspose2d(32, 3, kernel_size=4, stride=2,
padding=1), # 32 -> 64
        nn.Sigmoid()
    )

    def reparameterize(self, mu, logvar):
        if self.training:
            std = torch.exp(0.5 * logvar)
            eps = torch.randn_like(std)
            return mu + eps * std
        else:
            return mu

    def encode(self, x):
        x = self.encoder(x)
        x = x.view(x.size(0), -1)
        mu = self.fc_mu(x)
        logvar = self.fc_logvar(x)
        return mu, logvar

    def decode(self, z):
        x = self.decoder_input(z)
        x = x.view(x.size(0), 256, 4, 4)
        x = self.decoder(x)
        return x

    def forward(self, x):
        mu, logvar = self.encode(x)
        z = self.reparameterize(mu, logvar)
        return self.decode(z), mu, logvar

# Instantiate the model to verify structure
vae_model = VAE(latent_dim=128)
print(vae_model)

VAE(
  (encoder): Sequential(
    (0): Conv2d(3, 32, kernel_size=(4, 4), stride=(2, 2), padding=(1,
1))
    (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (2): LeakyReLU(negative_slope=0.2, inplace=True)
    (3): Conv2d(32, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1,

```

```

1))
    (4): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (5): LeakyReLU(negative_slope=0.2, inplace=True)
    (6): Conv2d(64, 128, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))
    (7): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (8): LeakyReLU(negative_slope=0.2, inplace=True)
    (9): Conv2d(128, 256, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))
    (10): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (11): LeakyReLU(negative_slope=0.2, inplace=True)
    )
    (fc_mu): Linear(in_features=4096, out_features=128, bias=True)
    (fc_logvar): Linear(in_features=4096, out_features=128, bias=True)
    (decoder_input): Linear(in_features=128, out_features=4096,
bias=True)
    (decoder): Sequential(
      (0): ConvTranspose2d(256, 128, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))
      (1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
      (2): ReLU(inplace=True)
      (3): ConvTranspose2d(128, 64, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))
      (4): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
      (5): ReLU(inplace=True)
      (6): ConvTranspose2d(64, 32, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))
      (7): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
      (8): ReLU(inplace=True)
      (9): ConvTranspose2d(32, 3, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))
      (10): Sigmoid()
    )
  )
)

```

Reasoning: Implementing the training workflow for the VAE model. This involves defining the loss function (combining reconstruction loss and KL divergence), initializing the model and optimizer, and running the training loop for 5 epochs using the previously created dataloader.

```

import torch.optim as optim

# Check for GPU
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")

```

```

print(f"Using device: {device}")

# Initialize model and optimizer
model = VAE(latent_dim=128).to(device)
optimizer = optim.Adam(model.parameters(), lr=1e-3)

# Define Loss Function
def loss_function(recon_x, x, mu, logvar):
    # Reconstruction loss (MSE)
    # reduction='sum' ensures it scales similarly to KLD
    MSE = F.mse_loss(recon_x, x, reduction='sum')

    # KL Divergence
    #  $0.5 * \sum(1 + \log(\sigma^2) - \mu^2 - \sigma^2)$ 
    KLD = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())

    return MSE + KLD

# Training Loop
epochs = 5
model.train()

print("Starting VAE training...")
for epoch in range(epochs):
    train_loss = 0
    for batch_idx, (data, _) in enumerate(dataloader):
        data = data.to(device)

        optimizer.zero_grad()

        recon_batch, mu, logvar = model(data)
        loss = loss_function(recon_batch, data, mu, logvar)

        loss.backward()
        train_loss += loss.item()
        optimizer.step()

        if batch_idx % 100 == 0:
            print(f'Epoch {epoch+1}
[batch_idx*len(data)]/[len(dataloader.dataset)] Loss:
{loss.item()/len(data):.4f}')

    print(f'====> Epoch: {epoch+1} Average loss: {train_loss /
len(dataloader.dataset):.4f}')

```

```

Using device: cuda
Starting VAE training...
Epoch 1 [0/162770] Loss: 1592.8649
Epoch 1 [6400/162770] Loss: 359.2393
Epoch 1 [12800/162770] Loss: 309.3353
Epoch 1 [19200/162770] Loss: 301.7740

```



```
Epoch 1 [25600/162770] Loss: 264.6439
Epoch 1 [32000/162770] Loss: 270.9055
Epoch 1 [38400/162770] Loss: 244.1530
Epoch 1 [44800/162770] Loss: 230.4097
Epoch 1 [51200/162770] Loss: 226.9161
Epoch 1 [57600/162770] Loss: 244.1601
Epoch 1 [64000/162770] Loss: 236.7159
Epoch 1 [70400/162770] Loss: 215.2279
Epoch 1 [76800/162770] Loss: 232.4265
Epoch 1 [83200/162770] Loss: 213.1579
Epoch 1 [89600/162770] Loss: 230.2442
Epoch 1 [96000/162770] Loss: 225.8381
Epoch 1 [102400/162770] Loss: 231.9338
Epoch 1 [108800/162770] Loss: 233.5865
Epoch 1 [115200/162770] Loss: 223.1697
Epoch 1 [121600/162770] Loss: 225.4326
Epoch 1 [128000/162770] Loss: 218.9953
Epoch 1 [134400/162770] Loss: 216.5358
Epoch 1 [140800/162770] Loss: 220.3188
Epoch 1 [147200/162770] Loss: 220.1432
Epoch 1 [153600/162770] Loss: 213.0195
Epoch 1 [160000/162770] Loss: 203.9446
====> Epoch: 1 Average loss: 247.7049
Epoch 2 [0/162770] Loss: 227.5977
Epoch 2 [6400/162770] Loss: 215.0998
Epoch 2 [12800/162770] Loss: 214.7168
Epoch 2 [19200/162770] Loss: 221.0088
Epoch 2 [25600/162770] Loss: 194.6798
Epoch 2 [32000/162770] Loss: 199.0235
Epoch 2 [38400/162770] Loss: 213.2731
Epoch 2 [44800/162770] Loss: 210.5651
Epoch 2 [51200/162770] Loss: 208.6142
Epoch 2 [57600/162770] Loss: 213.3364
Epoch 2 [64000/162770] Loss: 201.7514
Epoch 2 [70400/162770] Loss: 186.1395
Epoch 2 [76800/162770] Loss: 199.7092
Epoch 2 [83200/162770] Loss: 190.6600
Epoch 2 [89600/162770] Loss: 190.4880
Epoch 2 [96000/162770] Loss: 207.4631
Epoch 2 [102400/162770] Loss: 205.2829
Epoch 2 [108800/162770] Loss: 213.8417
Epoch 2 [115200/162770] Loss: 211.9508
Epoch 2 [121600/162770] Loss: 214.5592
Epoch 2 [128000/162770] Loss: 204.2356
Epoch 2 [134400/162770] Loss: 205.4363
Epoch 2 [140800/162770] Loss: 208.6780
Epoch 2 [147200/162770] Loss: 210.1507
Epoch 2 [153600/162770] Loss: 190.2908
Epoch 2 [160000/162770] Loss: 193.6252
```

```
====> Epoch: 2 Average loss: 205.6663
Epoch 3 [0/162770] Loss: 211.6532
Epoch 3 [6400/162770] Loss: 196.9123
Epoch 3 [12800/162770] Loss: 204.4644
Epoch 3 [19200/162770] Loss: 205.0269
Epoch 3 [25600/162770] Loss: 201.5197
Epoch 3 [32000/162770] Loss: 190.4644
Epoch 3 [38400/162770] Loss: 201.5741
Epoch 3 [44800/162770] Loss: 184.6531
Epoch 3 [51200/162770] Loss: 192.8651
Epoch 3 [57600/162770] Loss: 201.8669
Epoch 3 [64000/162770] Loss: 189.8502
Epoch 3 [70400/162770] Loss: 195.6109
Epoch 3 [76800/162770] Loss: 208.7139
Epoch 3 [83200/162770] Loss: 205.6473
Epoch 3 [89600/162770] Loss: 193.0598
Epoch 3 [96000/162770] Loss: 199.4350
Epoch 3 [102400/162770] Loss: 197.7137
Epoch 3 [108800/162770] Loss: 207.9081
Epoch 3 [115200/162770] Loss: 200.5110
Epoch 3 [121600/162770] Loss: 197.3123
Epoch 3 [128000/162770] Loss: 195.9280
Epoch 3 [134400/162770] Loss: 195.8942
Epoch 3 [140800/162770] Loss: 200.9460
Epoch 3 [147200/162770] Loss: 192.6833
Epoch 3 [153600/162770] Loss: 195.6440
Epoch 3 [160000/162770] Loss: 203.3846
====> Epoch: 3 Average loss: 199.2657
Epoch 4 [0/162770] Loss: 202.3535
Epoch 4 [6400/162770] Loss: 192.6177
Epoch 4 [12800/162770] Loss: 211.1486
Epoch 4 [19200/162770] Loss: 186.5634
Epoch 4 [25600/162770] Loss: 200.9262
Epoch 4 [32000/162770] Loss: 192.0333
Epoch 4 [38400/162770] Loss: 191.4011
Epoch 4 [44800/162770] Loss: 200.6657
Epoch 4 [51200/162770] Loss: 177.6590
Epoch 4 [57600/162770] Loss: 199.5869
Epoch 4 [64000/162770] Loss: 187.5701
Epoch 4 [70400/162770] Loss: 200.1452
Epoch 4 [76800/162770] Loss: 191.7482
Epoch 4 [83200/162770] Loss: 198.8639
Epoch 4 [89600/162770] Loss: 198.6143
Epoch 4 [96000/162770] Loss: 196.1615
Epoch 4 [102400/162770] Loss: 203.7421
Epoch 4 [108800/162770] Loss: 188.5216
Epoch 4 [115200/162770] Loss: 214.9430
Epoch 4 [121600/162770] Loss: 192.8132
Epoch 4 [128000/162770] Loss: 189.1935
```

```

Epoch 4 [134400/162770] Loss: 200.0648
Epoch 4 [140800/162770] Loss: 192.1985
Epoch 4 [147200/162770] Loss: 191.3095
Epoch 4 [153600/162770] Loss: 184.7988
Epoch 4 [160000/162770] Loss: 191.0537
====> Epoch: 4 Average loss: 195.5284
Epoch 5 [0/162770] Loss: 212.7182
Epoch 5 [6400/162770] Loss: 191.0568
Epoch 5 [12800/162770] Loss: 189.8726
Epoch 5 [19200/162770] Loss: 188.4250
Epoch 5 [25600/162770] Loss: 191.8532
Epoch 5 [32000/162770] Loss: 194.0622
Epoch 5 [38400/162770] Loss: 191.4741
Epoch 5 [44800/162770] Loss: 189.5387
Epoch 5 [51200/162770] Loss: 183.0632
Epoch 5 [57600/162770] Loss: 200.2629
Epoch 5 [64000/162770] Loss: 199.5594
Epoch 5 [70400/162770] Loss: 180.6450
Epoch 5 [76800/162770] Loss: 185.4145
Epoch 5 [83200/162770] Loss: 198.7004
Epoch 5 [89600/162770] Loss: 182.8950
Epoch 5 [96000/162770] Loss: 186.6487
Epoch 5 [102400/162770] Loss: 192.8550
Epoch 5 [108800/162770] Loss: 198.7262
Epoch 5 [115200/162770] Loss: 193.1144
Epoch 5 [121600/162770] Loss: 203.2660
Epoch 5 [128000/162770] Loss: 187.0650
Epoch 5 [134400/162770] Loss: 188.2140
Epoch 5 [140800/162770] Loss: 185.0892
Epoch 5 [147200/162770] Loss: 187.4808
Epoch 5 [153600/162770] Loss: 192.2664
Epoch 5 [160000/162770] Loss: 197.4912
====> Epoch: 5 Average loss: 192.5657

```

Reasoning: Implementing the visualization functions as requested. This involves visualising reconstructions of a random batch, performing linear interpolation in the latent space between two images, and calculating an attribute vector for 'Smiling' to modify a face image.

```

import matplotlib.pyplot as plt
import numpy as np

model.eval()

def show_image(img, title=""):
    img = img.cpu().numpy().transpose((1, 2, 0))
    plt.imshow(img)
    plt.axis("off")
    plt.title(title)

```

```

# 1. Visualize Reconstructions
print("Visualizing Reconstructions...")
with torch.no_grad():
    # Get a batch of images
    data_iter = iter(dataloader)
    images, labels = next(data_iter)
    images = images.to(device)

    # Reconstruct
    recon_images, _, _ = model(images)

    # Plot Original vs Reconstructed
    n = 8 # Number of images to display
    plt.figure(figsize=(16, 4))
    for i in range(n):
        # Original
        ax = plt.subplot(2, n, i + 1)
        show_image(images[i], title="Original")

        # Reconstructed
        ax = plt.subplot(2, n, i + 1 + n)
        show_image(recon_images[i], title="Reconstructed")
    plt.show()

# 2. Latent Interpolation
print("Visualizing Latent Interpolation...")
with torch.no_grad():
    z1, _ = model.encode(images[0].unsqueeze(0))
    z2, _ = model.encode(images[1].unsqueeze(0))

    steps = 10
    alphas = torch.linspace(0, 1, steps).to(device)
    interpolated_images = []

    for alpha in alphas:
        z_interp = (1 - alpha) * z1 + alpha * z2
        img_interp = model.decode(z_interp)
        interpolated_images.append(img_interp.squeeze(0))

    plt.figure(figsize=(20, 2))
    for i, img in enumerate(interpolated_images):
        ax = plt.subplot(1, steps, i + 1)
        show_image(img)
    plt.suptitle(f"Interpolation between Image 1 and Image 2")
    plt.show()

# 3. Attribute Modification (Smiling)
# 'Smiling' is typically index 31 in CelebA
smile_idx = 31

```

```

print("Calculating Attribute Vector for 'Smiling'...")
with torch.no_grad():
    z_pos = []
    z_neg = []
    limit_batches = 50 # Use a subset to calculate vector faster

    for i, (imgs, lbls) in enumerate(dataloader):
        if i >= limit_batches: break
        imgs = imgs.to(device)
        mu, _ = model.encode(imgs)

        # Separate based on attribute
        pos_mask = lbls[:, smile_idx] == 1
        neg_mask = lbls[:, smile_idx] == 0

        if pos_mask.sum() > 0:
            z_pos.append(mu[pos_mask])
        if neg_mask.sum() > 0:
            z_neg.append(mu[neg_mask])

    # Compute averages
    if z_pos and z_neg:
        z_pos = torch.cat(z_pos, dim=0).mean(dim=0)
        z_neg = torch.cat(z_neg, dim=0).mean(dim=0)
        smile_vector = z_pos - z_neg

    # Apply to a non-smiling face
    # Find a non-smiling face in current batch
    non_smile_indices = (labels[:, smile_idx] ==
0).nonzero(as_tuple=True)[0]

    if len(non_smile_indices) > 0:
        idx = non_smile_indices[0]
        original_img = images[idx]
        z_base, _ = model.encode(original_img.unsqueeze(0))

        # Add vector with a scaling factor
        scales = [-1.0, 0, 1.0, 2.0, 3.0]
        plt.figure(figsize=(15, 3))
        for i, s in enumerate(scales):
            z_mod = z_base + s * smile_vector
            img_mod = model.decode(z_mod)

            ax = plt.subplot(1, len(scales), i + 1)
            show_image(img_mod.squeeze(0), title=f"Scale {s}")
        plt.suptitle("Attribute Modification: Adding 'Smiling'
Vector")
        plt.show()
    else:

```

```
print("Could not calculate attribute vector (not enough samples).")
```

Visualizing Reconstructions...



Visualizing Latent Interpolation...



Calculating Attribute Vector for 'Smiling'...



Part B: Beta-VAE Implementation and Analysis

Subtask:

Explain the Beta-VAE objective, implement the training logic supporting different beta values, train a specific Beta-VAE model (beta=10) to demonstrate disentanglement, and visualize latent traversals.

1. Beta-VAE Concept and Objective

Beta-VAE (β -VAE) is a modification of the standard Variational Autoencoder framework that aims to discover disentangled latent factors. Disentanglement means that single latent units are sensitive to changes in single generative factors (like pose, skin color, or smile) while being relatively invariant to changes in other factors.

The objective function is a modified Evidence Lower Bound (ELBO) where the KL divergence term is weighted by a hyperparameter β :

$$L_{\beta\text{-VAE}} = E_{q_\phi(z \vee x)} [\log p_\theta(x \vee z)] - \beta \cdot D_{KL}$$

- **Reconstruction Term** (First term): Encourages high fidelity in reconstructing the input.
- **KL Divergence Term** (Second term): Constrains the capacity of the latent information channel.

The Role of β :

- When $\beta=1$, this corresponds to the standard VAE.
- When $\beta>1$, the model is forced to prioritize the KL divergence term, pushing the posterior $q(z \vee x)$ closer to the independent prior $p(z)$. This constraint limits the capacity of the latent bottleneck, encouraging the model to learn the most efficient data representation—often resulting in factorized (disentangled) latent dimensions.
- **Trade-off:** Higher β typically leads to better disentanglement but can result in blurrier reconstructions (loss of high-frequency details) because the model has less effective capacity to store information about the input.

Reasoning: I will implement the `train_beta_vae` function which modifies the loss calculation to include the beta weight on the KL divergence term. Then, I will train a Beta-VAE model with `beta=10` for 5 epochs to demonstrate the effect of high regularization.

```
def train_beta_vae(beta, epochs=5):
    print(f"\nStarting Beta-VAE training with beta={beta}...")

    # Initialize new model
    model_beta = VAE(latent_dim=128).to(device)
    optimizer = optim.Adam(model_beta.parameters(), lr=1e-3)
    model_beta.train()

    for epoch in range(epochs):
        train_loss = 0
        for batch_idx, (data, _) in enumerate(dataloader):
            data = data.to(device)
            optimizer.zero_grad()

            recon_batch, mu, logvar = model_beta(data)

            # Calculate loss components
            MSE = F.mse_loss(recon_batch, data, reduction='sum')
            KLD = -0.5 * torch.sum(1 + logvar - mu.pow(2) -
logvar.exp())

            # Beta-VAE Loss
            loss = MSE + beta * KLD

            loss.backward()
            train_loss += loss.item()
```

```

optimizer.step()

    if batch_idx % 100 == 0:
        # Print normalized loss for readability
        print(f'Epoch {epoch+1}
[batch_idx*len(data)]/[len(dataloader.dataset))] Loss:
{loss.item()/len(data):.4f}')

        print(f'====> Epoch: {epoch+1} Average loss: {train_loss /
len(dataloader.dataset):.4f}')
```

return model_beta

Train Beta-VAE with beta=10

beta_vae_model = train_beta_vae(beta=10, epochs=5)

Starting Beta-VAE training with beta=10...

```

Epoch 1 [0/162770] Loss: 1650.7899
Epoch 1 [6400/162770] Loss: 692.7953
Epoch 1 [12800/162770] Loss: 573.3315
Epoch 1 [19200/162770] Loss: 508.7409
Epoch 1 [25600/162770] Loss: 502.6106
Epoch 1 [32000/162770] Loss: 452.5489
Epoch 1 [38400/162770] Loss: 470.0809
Epoch 1 [44800/162770] Loss: 433.9919
Epoch 1 [51200/162770] Loss: 475.0965
Epoch 1 [57600/162770] Loss: 472.7202
Epoch 1 [64000/162770] Loss: 449.8520
Epoch 1 [70400/162770] Loss: 425.9923
Epoch 1 [76800/162770] Loss: 477.3568
Epoch 1 [83200/162770] Loss: 430.4725
Epoch 1 [89600/162770] Loss: 457.5925
Epoch 1 [96000/162770] Loss: 454.8442
Epoch 1 [102400/162770] Loss: 437.0110
Epoch 1 [108800/162770] Loss: 465.5670
Epoch 1 [115200/162770] Loss: 437.4462
Epoch 1 [121600/162770] Loss: 459.0110
Epoch 1 [128000/162770] Loss: 438.9694
Epoch 1 [134400/162770] Loss: 455.2005
Epoch 1 [140800/162770] Loss: 480.0364
Epoch 1 [147200/162770] Loss: 424.0531
Epoch 1 [153600/162770] Loss: 436.9267
Epoch 1 [160000/162770] Loss: 426.4895
====> Epoch: 1 Average loss: 482.6080
Epoch 2 [0/162770] Loss: 421.1071
Epoch 2 [6400/162770] Loss: 435.1567
Epoch 2 [12800/162770] Loss: 453.7712
Epoch 2 [19200/162770] Loss: 438.4396
Epoch 2 [25600/162770] Loss: 447.0362
```



```
Epoch 2 [32000/162770] Loss: 441.5255
Epoch 2 [38400/162770] Loss: 431.3105
Epoch 2 [44800/162770] Loss: 428.1633
Epoch 2 [51200/162770] Loss: 443.9423
Epoch 2 [57600/162770] Loss: 404.5013
Epoch 2 [64000/162770] Loss: 395.6635
Epoch 2 [70400/162770] Loss: 429.0514
Epoch 2 [76800/162770] Loss: 447.0113
Epoch 2 [83200/162770] Loss: 422.6242
Epoch 2 [89600/162770] Loss: 407.7366
Epoch 2 [96000/162770] Loss: 440.4193
Epoch 2 [102400/162770] Loss: 432.8201
Epoch 2 [108800/162770] Loss: 432.3341
Epoch 2 [115200/162770] Loss: 378.3485
Epoch 2 [121600/162770] Loss: 424.2733
Epoch 2 [128000/162770] Loss: 421.8125
Epoch 2 [134400/162770] Loss: 413.6191
Epoch 2 [140800/162770] Loss: 392.1958
Epoch 2 [147200/162770] Loss: 416.4655
Epoch 2 [153600/162770] Loss: 406.7837
Epoch 2 [160000/162770] Loss: 415.3506
====> Epoch: 2 Average loss: 430.8898
Epoch 3 [0/162770] Loss: 419.1321
Epoch 3 [6400/162770] Loss: 455.7249
Epoch 3 [12800/162770] Loss: 422.0023
Epoch 3 [19200/162770] Loss: 433.3774
Epoch 3 [25600/162770] Loss: 432.9113
Epoch 3 [32000/162770] Loss: 441.1545
Epoch 3 [38400/162770] Loss: 432.2941
Epoch 3 [44800/162770] Loss: 433.1760
Epoch 3 [51200/162770] Loss: 437.8911
Epoch 3 [57600/162770] Loss: 415.3226
Epoch 3 [64000/162770] Loss: 436.8010
Epoch 3 [70400/162770] Loss: 427.4412
Epoch 3 [76800/162770] Loss: 448.9692
Epoch 3 [83200/162770] Loss: 410.7578
Epoch 3 [89600/162770] Loss: 425.1799
Epoch 3 [96000/162770] Loss: 433.5683
Epoch 3 [102400/162770] Loss: 401.1943
Epoch 3 [108800/162770] Loss: 419.1085
Epoch 3 [115200/162770] Loss: 450.2666
Epoch 3 [121600/162770] Loss: 423.7939
Epoch 3 [128000/162770] Loss: 433.3821
Epoch 3 [134400/162770] Loss: 390.4178
Epoch 3 [140800/162770] Loss: 411.8977
Epoch 3 [147200/162770] Loss: 443.5841
Epoch 3 [153600/162770] Loss: 423.9115
Epoch 3 [160000/162770] Loss: 426.2889
====> Epoch: 3 Average loss: 424.0901
```

```
Epoch 4 [0/162770] Loss: 403.8444
Epoch 4 [6400/162770] Loss: 392.7150
Epoch 4 [12800/162770] Loss: 432.2758
Epoch 4 [19200/162770] Loss: 462.6322
Epoch 4 [25600/162770] Loss: 433.8685
Epoch 4 [32000/162770] Loss: 421.2604
Epoch 4 [38400/162770] Loss: 443.4559
Epoch 4 [44800/162770] Loss: 414.7825
Epoch 4 [51200/162770] Loss: 419.8247
Epoch 4 [57600/162770] Loss: 444.2333
Epoch 4 [64000/162770] Loss: 415.8748
Epoch 4 [70400/162770] Loss: 394.9742
Epoch 4 [76800/162770] Loss: 415.4728
Epoch 4 [83200/162770] Loss: 386.0418
Epoch 4 [89600/162770] Loss: 411.5795
Epoch 4 [96000/162770] Loss: 425.3138
Epoch 4 [102400/162770] Loss: 404.5393
Epoch 4 [108800/162770] Loss: 434.1639
Epoch 4 [115200/162770] Loss: 429.0736
Epoch 4 [121600/162770] Loss: 428.8761
Epoch 4 [128000/162770] Loss: 386.4613
Epoch 4 [134400/162770] Loss: 439.4922
Epoch 4 [140800/162770] Loss: 441.6085
Epoch 4 [147200/162770] Loss: 395.9907
Epoch 4 [153600/162770] Loss: 408.8757
Epoch 4 [160000/162770] Loss: 441.4558
====> Epoch: 4 Average loss: 420.8229
Epoch 5 [0/162770] Loss: 437.0216
Epoch 5 [6400/162770] Loss: 458.5834
Epoch 5 [12800/162770] Loss: 414.5124
Epoch 5 [19200/162770] Loss: 406.4426
Epoch 5 [25600/162770] Loss: 400.4283
Epoch 5 [32000/162770] Loss: 416.2152
Epoch 5 [38400/162770] Loss: 447.6714
Epoch 5 [44800/162770] Loss: 412.8635
Epoch 5 [51200/162770] Loss: 408.3680
Epoch 5 [57600/162770] Loss: 417.2931
Epoch 5 [64000/162770] Loss: 427.8489
Epoch 5 [70400/162770] Loss: 433.8430
Epoch 5 [76800/162770] Loss: 422.4915
Epoch 5 [83200/162770] Loss: 437.6226
Epoch 5 [89600/162770] Loss: 404.4666
Epoch 5 [96000/162770] Loss: 454.8136
Epoch 5 [102400/162770] Loss: 421.0783
Epoch 5 [108800/162770] Loss: 389.7955
Epoch 5 [115200/162770] Loss: 422.5309
Epoch 5 [121600/162770] Loss: 433.0762
Epoch 5 [128000/162770] Loss: 425.7053
Epoch 5 [134400/162770] Loss: 434.2866
```

```
Epoch 5 [140800/162770] Loss: 418.8761
Epoch 5 [147200/162770] Loss: 415.1882
Epoch 5 [153600/162770] Loss: 400.0200
Epoch 5 [160000/162770] Loss: 413.6938
====> Epoch: 5 Average loss: 418.3295
```

Reasoning: Implementing the `visualize_latent_traversal` function to demonstrate the disentanglement properties of the Beta-VAE. This function will vary specific latent dimensions while keeping others fixed, allowing us to observe if specific dimensions control distinct facial features. I will then apply this function to a sample image using the trained Beta-VAE model.

Task

Visualize the disentanglement properties of the trained Beta-VAE model by implementing a latent traversal function. This function should:

1. Identify the most "active" latent dimensions (those with the highest variance across a batch) to avoid visualizing collapsed dimensions.
2. Take a sample face image, encode it, and then traverse (vary) specific latent dimensions one by one over a range (e.g., $[-3, 3]$).
3. Decode these modified latent vectors and plot the resulting images in a grid to show how specific dimensions control distinct facial features.

```
import torch
import matplotlib.pyplot as plt
import numpy as np

def get_active_dimensions(model, dataloader, limit=10):
    """Finds latent dimensions with highest variance across a
    batch."""
    model.eval()
    with torch.no_grad():
        images, _ = next(iter(dataloader))
        images = images.to(device)
        mu, _ = model.encode(images)
        # Calculate variance across the batch for each dimension
        vars = torch.var(mu, dim=0)
        # Get indices of dimensions with highest variance
        _, sorted_indices = torch.sort(vars, descending=True)
        return sorted_indices[:limit].cpu().numpy()

def visualize_latent_traversal(model, image, active_dims,
    title_prefix=""):
    """Varies latent dimensions to visualize disentanglement."""
    model.eval()
    with torch.no_grad():
        # Encode the single image
        img_tensor = image.unsqueeze(0).to(device)
```

```

z_orig, _ = model.encode(img_tensor)

# Define range for traversal
traversal_range = torch.linspace(-3, 3, 10).to(device)

plt.figure(figsize=(15, len(active_dims) * 1.5))

for row_idx, dim_idx in enumerate(active_dims):
    # Create a batch of latent vectors for this dimension's
traversal
    z_traversal = z_orig.repeat(len(traversal_range), 1)

    # Vary the specific dimension
    # Note: We replace the value, seeing how the decoder
responds
    z_traversal[:, dim_idx] = traversal_range

    # Decode
    recon_traversal = model.decode(z_traversal)

    # Plot
    for col_idx in range(len(traversal_range)):
        ax = plt.subplot(len(active_dims),
len(traversal_range), row_idx * len(traversal_range) + col_idx + 1)
        img =
recon_traversal[col_idx].cpu().numpy().transpose((1, 2, 0))
        plt.imshow(img)
        plt.axis("off")
        if col_idx == 0:
            plt.text(-0.5, 0.5, f"Dim {dim_idx}", fontsize=12,
transform=ax.transAxes, va='center', ha='right')

    plt.suptitle(f"{title_prefix} Latent Traversal on Top
{len(active_dims)} Active Dimensions", y=1.02)
    plt.tight_layout()
    plt.show()

# Execute Visualization
print("Identifying active dimensions...")
active_dims = get_active_dimensions(beta_vae_model, dataloader,
limit=8)
print(f"Active Dimensions: {active_dims}")

# Get a sample image
data_iter = iter(dataloader)
images, _ = next(data_iter)
sample_image = images[0]

print("Visualizing Traversal...")

```

```
visualize_latent_traversal(beta_vae_model, sample_image, active_dims,
title_prefix="Beta-VAE (beta=10)")
```

Visualize Beta-VAE Disentanglement

Subtask:

Implement and execute latent traversal functions to visualize how specific latent dimensions affect the generated images in the Beta-VAE model, demonstrating disentanglement.

Reasoning: Implementing the `get_active_dimensions` and `visualize_latent_traversal` functions to identify and visualize the most active latent dimensions of the trained Beta-VAE model, effectively demonstrating the learned disentanglement.

```
import matplotlib.pyplot as plt
import numpy as np
import torch

def get_active_dimensions(model, dataloader, limit=10):
    """Identify latent dimensions with highest variance across a
    batch."""
    model.eval()
    with torch.no_grad():
        # Get a batch of data
        data_iter = iter(dataloader)
        images, _ = next(data_iter)
        images = images.to(device)

        # Encode to get mu
        mu, _ = model.encode(images)

        # Calculate variance across the batch for each latent
        dimension
        # Shape: (batch_size, latent_dim) -> (latent_dim)
        variances = torch.var(mu, dim=0)

        # Sort and get top indices
        _, sorted_indices = torch.sort(variances, descending=True)
        active_dims = sorted_indices[:limit].tolist()

    return active_dims

def visualize_latent_traversal(model, image, active_dims,
title_prefix=""):
    """Vary specific dimensions and plot results."""
    model.eval()
    with torch.no_grad():
        image = image.to(device).unsqueeze(0)
```

```

z_base, _ = model.encode(image)

# Traversal range
traversal_range = torch.linspace(-3, 3, 10).to(device)

plt.figure(figsize=(15, 1.5 * len(active_dims)))

for i, dim_idx in enumerate(active_dims):
    traversal_vectors = []
    for val in traversal_range:
        z_new = z_base.clone()
        z_new[0, dim_idx] = val
        traversal_vectors.append(z_new)

    traversal_vectors = torch.cat(traversal_vectors, dim=0)
    decoded_images = model.decode(traversal_vectors)

    # Plotting
    for j in range(len(traversal_range)):
        ax = plt.subplot(len(active_dims), 10, i * 10 + j + 1)
        img = decoded_images[j].cpu().permute(1, 2, 0).numpy()
        plt.imshow(img)
        plt.axis('off')
        if j == 0:
            ax.set_title(f'Dim {dim_idx}', fontsize=10,
loc='left')

    plt.suptitle(f"{title_prefix} Latent Traversal on Top
{len(active_dims)} Active Dimensions")
    plt.tight_layout()
    plt.show()

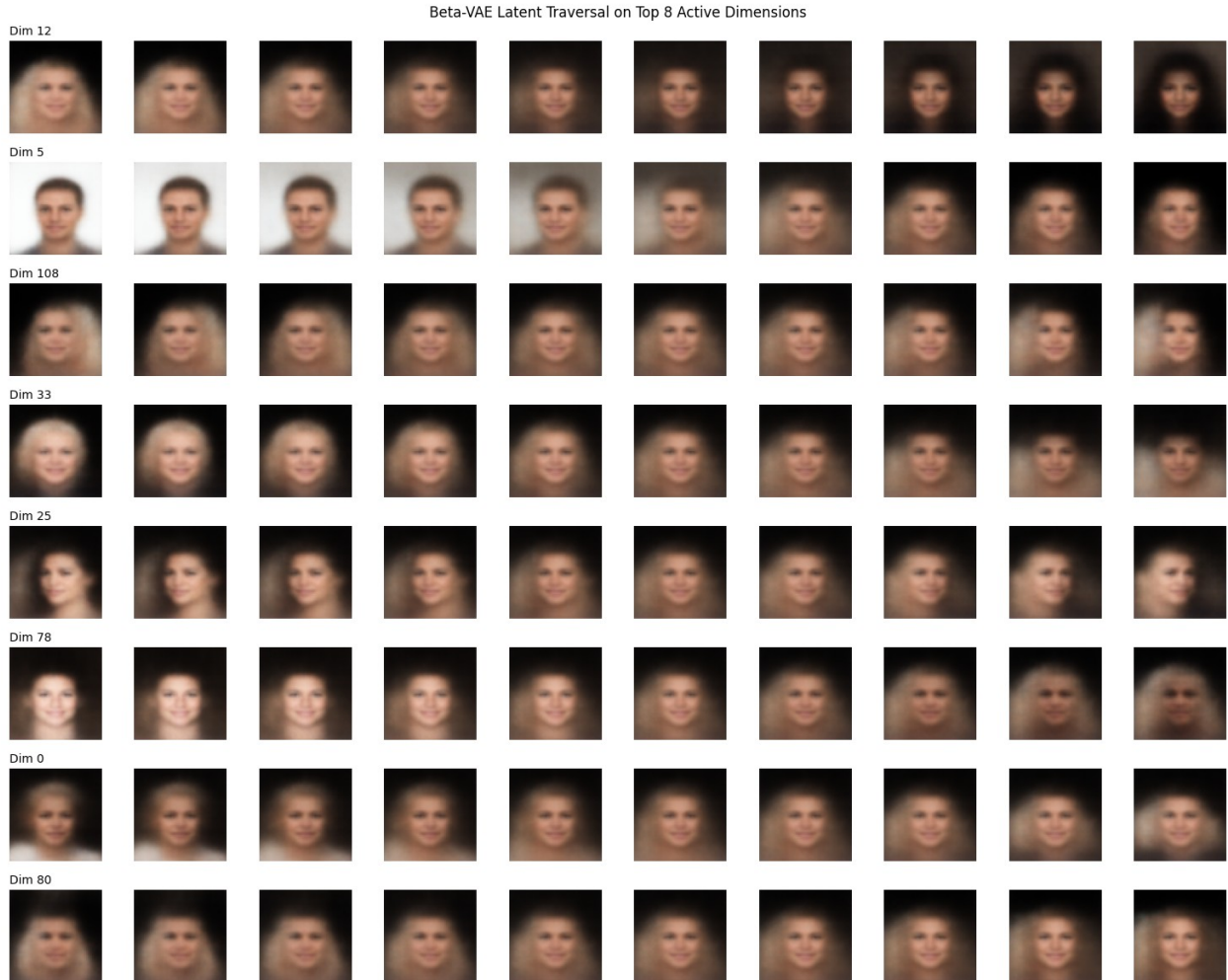
# Execute Analysis
print("Identifying active dimensions...")
top_dims = get_active_dimensions(beta_vae_model, dataloader, limit=8)
print(f"Top active dimensions: {top_dims}")

# Get a sample image
img_sample = next(iter(dataloader))[0][0]

# Visualize
print("Visualizing latent traversal...")
visualize_latent_traversal(beta_vae_model, img_sample, top_dims,
title_prefix="Beta-VAE")

Identifying active dimensions...
Top active dimensions: [12, 5, 108, 33, 25, 78, 0, 80]
Visualizing latent traversal...

```



Part C: VQ-VAE Implementation

Subtask:

Derive the VQ-VAE loss function, implement the Vector Quantizer with the straight-through estimator, build the VQ-VAE architecture, and train it on CelebA.

1. VQ-VAE Objective Function

The Vector Quantized-Variational Autoencoder (VQ-VAE) differs from standard VAEs by using a discrete latent representation. The objective function is given by:

$$L = \log p(x | z_q) + \frac{1}{2} \| \text{sg}(z_e) - e \|_2^2 + \beta \frac{1}{2} \| z_e - \text{sg}(e) \|_2^2$$

Where:

- $z_e(x)$ is the output of the encoder (continuous).
- e is the nearest embedding vector in the codebook.
- $z_q(x)$ is the quantized latent vector (the chosen embedding).

- $\text{sg}[\cdot]$ stands for the **stop-gradient** operator.

Terms:

1. **Reconstruction Loss:** Ensures the decoder can reconstruct the input x from the quantized latent z_q . Typically essentially Mean Squared Error (MSE).
2. **Codebook Loss:** Updates the codebook vectors (embeddings) to move closer to the encoder outputs z_e . Since gradients don't flow through the discrete selection, we use the stop-gradient on z_e .
3. **Commitment Loss:** Constrains the encoder output z_e to stay close to the chosen embedding e , preventing it from growing arbitrarily large or fluctuating wildly. β is a hyperparameter (usually 0.25).

Straight-Through Estimator: Since the mapping from z_e to z_q (argmin) is non-differentiable, we approximate the gradient during backpropagation by copying the gradients from z_q to z_e :

$$\frac{\partial L}{\partial z_e} \approx \frac{\partial L}{\partial z_q}$$

In PyTorch, this is implemented as `z_q = z_e + (z_q - z_e).detach()`.

Reasoning: Implementing the `VectorQuantizer` and `VQVAE` classes. The `VectorQuantizer` handles the discrete bottleneck, codebook updates, and the straight-through estimator. The `VQVAE` combines the encoder, quantizer, and decoder.

```
class VectorQuantizer(nn.Module):
    def __init__(self, num_embeddings, embedding_dim,
                 commitment_cost=0.25):
        super(VectorQuantizer, self).__init__()
        self._embedding_dim = embedding_dim
        self._num_embeddings = num_embeddings
        self._commitment_cost = commitment_cost

        self._embedding = nn.Embedding(self._num_embeddings,
                                       self._embedding_dim)
        self._embedding.weight.data.uniform_(-1/self._num_embeddings,
                                             1/self._num_embeddings)

    def forward(self, inputs):
        # convert inputs from BCHW -> BHCW
        inputs = inputs.permute(0, 2, 3, 1).contiguous()
        input_shape = inputs.shape

        # Flatten input
        flat_input = inputs.view(-1, self._embedding_dim)

        # Calculate distances
        distances = (torch.sum(flat_input**2, dim=1, keepdim=True)
                     + torch.sum(self._embedding.weight**2, dim=1)
                     - 2 * torch.matmul(flat_input,
```



```

self._embedding.weight.t()))

    # Encoding
    encoding_indices = torch.argmax(distances, dim=1).unsqueeze(1)
    encodings = torch.zeros(encoding_indices.shape[0],
self._num_embeddings, device=inputs.device)
    encodings.scatter_(1, encoding_indices, 1)

    # Quantize and unflatten
    quantized = torch.matmul(encodings,
self._embedding.weight).view(input_shape)

    # Loss
    e_latent_loss = F.mse_loss(quantized.detach(), inputs)
    q_latent_loss = F.mse_loss(quantized, inputs.detach())
    loss = q_latent_loss + self._commitment_cost * e_latent_loss

    # Straight Through Estimator
    quantized = inputs + (quantized - inputs).detach()
    avg_probs = torch.mean(encodings, dim=0)
    perplexity = torch.exp(-torch.sum(avg_probs *
torch.log(avg_probs + 1e-10)))

    # convert quantized from BHWC -> BCHW
    return loss, quantized.permute(0, 3, 1, 2).contiguous(),
perplexity, encoding_indices

class VQVAE(nn.Module):
    def __init__(self, num_hiddens, num_residual_hiddens,
num_embeddings, embedding_dim, commitment_cost=0.25):
        super(VQVAE, self).__init__()

        # Encoder
        self._encoder = nn.Sequential(
            nn.Conv2d(3, num_hiddens // 2, kernel_size=4, stride=2,
padding=1),
            nn.ReLU(True),
            nn.Conv2d(num_hiddens // 2, num_hiddens, kernel_size=4,
stride=2, padding=1),
            nn.ReLU(True),
            nn.Conv2d(num_hiddens, num_hiddens, kernel_size=3,
stride=1, padding=1),
            nn.ReLU(True),
            nn.Conv2d(num_hiddens, embedding_dim, kernel_size=1,
stride=1)
        )

        self._pre_vq_conv = nn.Conv2d(embedding_dim, embedding_dim, 1)
        self._vq_vae = VectorQuantizer(num_embeddings, embedding_dim,
commitment_cost)

```

```

    # Decoder
    self._decoder = nn.Sequential(
        nn.Conv2d(embedding_dim, num_hiddens, kernel_size=3,
stride=1, padding=1),
        nn.ReLU(True),
        nn.ConvTranspose2d(num_hiddens, num_hiddens // 2,
kernel_size=4, stride=2, padding=1),
        nn.ReLU(True),
        nn.ConvTranspose2d(num_hiddens // 2, 3, kernel_size=4,
stride=2, padding=1),
        nn.Sigmoid() # Normalize to [0, 1]
    )

    def forward(self, x):
        z = self._encoder(x)
        z = self._pre_vq_conv(z)
        loss, quantized, perplexity, _ = self._vq_vae(z)
        x_recon = self._decoder(quantized)
        return loss, x_recon, perplexity

# Initialize Model
# K=512, D=64 as requested
vqvae_model = VQVAE(num_hiddens=128, num_residual_hiddens=32,
num_embeddings=512, embedding_dim=64).to(device)
print(vqvae_model)

VQVAE(
  (_encoder): Sequential(
    (0): Conv2d(3, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1,
1))
    (1): ReLU(inplace=True)
    (2): Conv2d(64, 128, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))
    (3): ReLU(inplace=True)
    (4): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1),
padding=(1, 1))
    (5): ReLU(inplace=True)
    (6): Conv2d(128, 64, kernel_size=(1, 1), stride=(1, 1))
  )
  (_pre_vq_conv): Conv2d(64, 64, kernel_size=(1, 1), stride=(1, 1))
  (_vq_vae): VectorQuantizer(
    (_embedding): Embedding(512, 64)
  )
  (_decoder): Sequential(
    (0): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1),
padding=(1, 1))
    (1): ReLU(inplace=True)
    (2): ConvTranspose2d(128, 64, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))

```

```

        (3): ReLU(inplace=True)
        (4): ConvTranspose2d(64, 3, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1))
        (5): Sigmoid()
    )
)

```

Reasoning: Implementing the training loop for the VQ-VAE model. This involves defining the optimizer, calculating the total loss (Reconstruction Loss + VQ Loss), and iterating through the dataset for 5 epochs.

```

import torch.optim as optim

# Initialize optimizer
optimizer = optim.Adam(vqvae_model.parameters(), lr=1e-3)

print("Starting VQ-VAE training...")
num_training_updates = 5 * len(dataloader)

# Train for 5 epochs
for epoch in range(5):
    train_loss = 0
    train_recon_loss = 0
    train_vq_loss = 0

    for batch_idx, (data, _) in enumerate(dataloader):
        data = data.to(device)
        optimizer.zero_grad()

        # Forward pass
        vq_loss, data_recon, perplexity = vqvae_model(data)

        # Reconstruction Loss
        recon_loss = F.mse_loss(data_recon, data)

        # Total Loss
        loss = recon_loss + vq_loss

        loss.backward()
        optimizer.step()

        train_loss += loss.item()
        train_recon_loss += recon_loss.item()
        train_vq_loss += vq_loss.item()

        if batch_idx % 100 == 0:
            print(f'Epoch {epoch+1}
[{batch_idx*len(data)}/{len(dataloader.dataset)}] '
                  f'Loss: {loss.item():.4f} | Recon:
{recon_loss.item():.4f} | VQ: {vq_loss.item():.4f}')

```

```

avg_loss = train_loss / len(dataloader)
avg_recon = train_recon_loss / len(dataloader)
avg_vq = train_vq_loss / len(dataloader)

print(f'====> Epoch: {epoch+1} Average Loss: {avg_loss:.4f}
(Recon: {avg_recon:.4f}, VQ: {avg_vq:.4f})')

```

Starting VQ-VAE training...

Epoch	Progress	Loss	Recon	VQ
Epoch 1	[0/162770]	Loss: 0.0977	Recon: 0.0908	VQ: 0.0069
Epoch 1	[6400/162770]	Loss: 0.3722	Recon: 0.0227	VQ: 0.3494
Epoch 1	[12800/162770]	Loss: 0.0492	Recon: 0.0145	VQ: 0.0347
Epoch 1	[19200/162770]	Loss: 0.0134	Recon: 0.0092	VQ: 0.0042
Epoch 1	[25600/162770]	Loss: 0.0118	Recon: 0.0074	VQ: 0.0045
Epoch 1	[32000/162770]	Loss: 0.0090	Recon: 0.0062	VQ: 0.0027
Epoch 1	[38400/162770]	Loss: 0.0079	Recon: 0.0054	VQ: 0.0026
Epoch 1	[44800/162770]	Loss: 0.0078	Recon: 0.0052	VQ: 0.0026
Epoch 1	[51200/162770]	Loss: 0.0078	Recon: 0.0050	VQ: 0.0028
Epoch 1	[57600/162770]	Loss: 0.0075	Recon: 0.0048	VQ: 0.0027
Epoch 1	[64000/162770]	Loss: 0.0068	Recon: 0.0044	VQ: 0.0024
Epoch 1	[70400/162770]	Loss: 0.0062	Recon: 0.0038	VQ: 0.0023
Epoch 1	[76800/162770]	Loss: 0.0076	Recon: 0.0045	VQ: 0.0031
Epoch 1	[83200/162770]	Loss: 0.0070	Recon: 0.0038	VQ: 0.0032
Epoch 1	[89600/162770]	Loss: 0.0061	Recon: 0.0033	VQ: 0.0028
Epoch 1	[96000/162770]	Loss: 0.0064	Recon: 0.0035	VQ: 0.0030
Epoch 1	[102400/162770]	Loss: 0.0062	Recon: 0.0032	VQ: 0.0030
Epoch 1	[108800/162770]	Loss: 0.0064	Recon: 0.0035	VQ: 0.0029
Epoch 1	[115200/162770]	Loss: 0.0066	Recon: 0.0037	VQ: 0.0030
Epoch 1	[121600/162770]	Loss: 0.0063	Recon: 0.0034	VQ: 0.0029
Epoch 1	[128000/162770]	Loss: 0.0064	Recon: 0.0033	VQ: 0.0031
Epoch 1	[134400/162770]	Loss: 0.0058	Recon: 0.0032	VQ: 0.0025
Epoch 1	[140800/162770]	Loss: 0.0058	Recon: 0.0032	VQ: 0.0027
Epoch 1	[147200/162770]	Loss: 0.0061	Recon: 0.0032	VQ: 0.0029
Epoch 1	[153600/162770]	Loss: 0.0058	Recon: 0.0032	VQ: 0.0026
Epoch 1	[160000/162770]	Loss: 0.0057	Recon: 0.0032	VQ: 0.0025
====> Epoch: 1 Average Loss: 0.0401 (Recon: 0.0067, VQ: 0.0334)				
Epoch 2	[0/162770]	Loss: 0.0053	Recon: 0.0030	VQ: 0.0022
Epoch 2	[6400/162770]	Loss: 0.0061	Recon: 0.0033	VQ: 0.0028
Epoch 2	[12800/162770]	Loss: 0.0050	Recon: 0.0027	VQ: 0.0022
Epoch 2	[19200/162770]	Loss: 0.0050	Recon: 0.0029	VQ: 0.0022
Epoch 2	[25600/162770]	Loss: 0.0051	Recon: 0.0028	VQ: 0.0023
Epoch 2	[32000/162770]	Loss: 0.0050	Recon: 0.0029	VQ: 0.0021
Epoch 2	[38400/162770]	Loss: 0.0051	Recon: 0.0029	VQ: 0.0022
Epoch 2	[44800/162770]	Loss: 0.0047	Recon: 0.0027	VQ: 0.0020
Epoch 2	[51200/162770]	Loss: 0.0050	Recon: 0.0029	VQ: 0.0021
Epoch 2	[57600/162770]	Loss: 0.0050	Recon: 0.0030	VQ: 0.0020
Epoch 2	[64000/162770]	Loss: 0.0051	Recon: 0.0032	VQ: 0.0020
Epoch 2	[70400/162770]	Loss: 0.0048	Recon: 0.0028	VQ: 0.0020
Epoch 2	[76800/162770]	Loss: 0.0051	Recon: 0.0029	VQ: 0.0022
Epoch 2	[83200/162770]	Loss: 0.0048	Recon: 0.0028	VQ: 0.0020

Epoch 2	[89600/162770]	Loss: 0.0048	Recon: 0.0027	VQ: 0.0021
Epoch 2	[96000/162770]	Loss: 0.0051	Recon: 0.0026	VQ: 0.0025
Epoch 2	[102400/162770]	Loss: 0.0056	Recon: 0.0029	VQ: 0.0027
Epoch 2	[108800/162770]	Loss: 0.0056	Recon: 0.0025	VQ: 0.0031
Epoch 2	[115200/162770]	Loss: 0.0055	Recon: 0.0026	VQ: 0.0028
Epoch 2	[121600/162770]	Loss: 0.0063	Recon: 0.0027	VQ: 0.0035
Epoch 2	[128000/162770]	Loss: 0.0059	Recon: 0.0026	VQ: 0.0032
Epoch 2	[134400/162770]	Loss: 0.0056	Recon: 0.0024	VQ: 0.0032
Epoch 2	[140800/162770]	Loss: 0.0058	Recon: 0.0026	VQ: 0.0032
Epoch 2	[147200/162770]	Loss: 0.0057	Recon: 0.0026	VQ: 0.0032
Epoch 2	[153600/162770]	Loss: 0.0058	Recon: 0.0024	VQ: 0.0034
Epoch 2	[160000/162770]	Loss: 0.0054	Recon: 0.0023	VQ: 0.0031
====> Epoch: 2 Average Loss: 0.0054 (Recon: 0.0028, VQ: 0.0026)				
Epoch 3	[0/162770]	Loss: 0.0057	Recon: 0.0023	VQ: 0.0034
Epoch 3	[6400/162770]	Loss: 0.0053	Recon: 0.0023	VQ: 0.0031
Epoch 3	[12800/162770]	Loss: 0.0055	Recon: 0.0023	VQ: 0.0032
Epoch 3	[19200/162770]	Loss: 0.0060	Recon: 0.0025	VQ: 0.0035
Epoch 3	[25600/162770]	Loss: 0.0057	Recon: 0.0023	VQ: 0.0034
Epoch 3	[32000/162770]	Loss: 0.0051	Recon: 0.0022	VQ: 0.0029
Epoch 3	[38400/162770]	Loss: 0.0053	Recon: 0.0022	VQ: 0.0031
Epoch 3	[44800/162770]	Loss: 0.0055	Recon: 0.0022	VQ: 0.0033
Epoch 3	[51200/162770]	Loss: 0.0057	Recon: 0.0024	VQ: 0.0033
Epoch 3	[57600/162770]	Loss: 0.0054	Recon: 0.0022	VQ: 0.0033
Epoch 3	[64000/162770]	Loss: 0.0055	Recon: 0.0024	VQ: 0.0031
Epoch 3	[70400/162770]	Loss: 0.0054	Recon: 0.0024	VQ: 0.0031
Epoch 3	[76800/162770]	Loss: 0.0059	Recon: 0.0022	VQ: 0.0037
Epoch 3	[83200/162770]	Loss: 0.0056	Recon: 0.0023	VQ: 0.0033
Epoch 3	[89600/162770]	Loss: 0.0049	Recon: 0.0020	VQ: 0.0029
Epoch 3	[96000/162770]	Loss: 0.0051	Recon: 0.0022	VQ: 0.0029
Epoch 3	[102400/162770]	Loss: 0.0048	Recon: 0.0019	VQ: 0.0028
Epoch 3	[108800/162770]	Loss: 0.0049	Recon: 0.0020	VQ: 0.0028
Epoch 3	[115200/162770]	Loss: 0.0048	Recon: 0.0020	VQ: 0.0028
Epoch 3	[121600/162770]	Loss: 0.0052	Recon: 0.0021	VQ: 0.0031
Epoch 3	[128000/162770]	Loss: 0.0050	Recon: 0.0021	VQ: 0.0029
Epoch 3	[134400/162770]	Loss: 0.0046	Recon: 0.0019	VQ: 0.0027
Epoch 3	[140800/162770]	Loss: 0.0053	Recon: 0.0022	VQ: 0.0031
Epoch 3	[147200/162770]	Loss: 0.0052	Recon: 0.0022	VQ: 0.0030
Epoch 3	[153600/162770]	Loss: 0.0049	Recon: 0.0021	VQ: 0.0028
Epoch 3	[160000/162770]	Loss: 0.0054	Recon: 0.0023	VQ: 0.0031
====> Epoch: 3 Average Loss: 0.0052 (Recon: 0.0022, VQ: 0.0031)				
Epoch 4	[0/162770]	Loss: 0.0055	Recon: 0.0022	VQ: 0.0033
Epoch 4	[6400/162770]	Loss: 0.0051	Recon: 0.0022	VQ: 0.0028
Epoch 4	[12800/162770]	Loss: 0.0046	Recon: 0.0020	VQ: 0.0026
Epoch 4	[19200/162770]	Loss: 0.0048	Recon: 0.0020	VQ: 0.0028
Epoch 4	[25600/162770]	Loss: 0.0049	Recon: 0.0022	VQ: 0.0027
Epoch 4	[32000/162770]	Loss: 0.0054	Recon: 0.0024	VQ: 0.0029
Epoch 4	[38400/162770]	Loss: 0.0047	Recon: 0.0020	VQ: 0.0026
Epoch 4	[44800/162770]	Loss: 0.0044	Recon: 0.0019	VQ: 0.0025
Epoch 4	[51200/162770]	Loss: 0.0045	Recon: 0.0020	VQ: 0.0025

```

Epoch 4 [57600/162770] Loss: 0.0050 | Recon: 0.0021 | VQ: 0.0029
Epoch 4 [64000/162770] Loss: 0.0043 | Recon: 0.0017 | VQ: 0.0025
Epoch 4 [70400/162770] Loss: 0.0048 | Recon: 0.0021 | VQ: 0.0027
Epoch 4 [76800/162770] Loss: 0.0042 | Recon: 0.0018 | VQ: 0.0024
Epoch 4 [83200/162770] Loss: 0.0050 | Recon: 0.0023 | VQ: 0.0028
Epoch 4 [89600/162770] Loss: 0.0041 | Recon: 0.0017 | VQ: 0.0024
Epoch 4 [96000/162770] Loss: 0.0044 | Recon: 0.0019 | VQ: 0.0025
Epoch 4 [102400/162770] Loss: 0.0043 | Recon: 0.0018 | VQ: 0.0024
Epoch 4 [108800/162770] Loss: 0.0042 | Recon: 0.0018 | VQ: 0.0024
Epoch 4 [115200/162770] Loss: 0.0042 | Recon: 0.0018 | VQ: 0.0023
Epoch 4 [121600/162770] Loss: 0.0044 | Recon: 0.0020 | VQ: 0.0024
Epoch 4 [128000/162770] Loss: 0.0046 | Recon: 0.0020 | VQ: 0.0026
Epoch 4 [134400/162770] Loss: 0.0044 | Recon: 0.0019 | VQ: 0.0024
Epoch 4 [140800/162770] Loss: 0.0048 | Recon: 0.0021 | VQ: 0.0027
Epoch 4 [147200/162770] Loss: 0.0043 | Recon: 0.0019 | VQ: 0.0025
Epoch 4 [153600/162770] Loss: 0.0042 | Recon: 0.0019 | VQ: 0.0023
Epoch 4 [160000/162770] Loss: 0.0042 | Recon: 0.0018 | VQ: 0.0024
====> Epoch: 4 Average Loss: 0.0046 (Recon: 0.0020, VQ: 0.0026)
Epoch 5 [0/162770] Loss: 0.0042 | Recon: 0.0018 | VQ: 0.0023
Epoch 5 [6400/162770] Loss: 0.0044 | Recon: 0.0019 | VQ: 0.0025
Epoch 5 [12800/162770] Loss: 0.0045 | Recon: 0.0020 | VQ: 0.0025
Epoch 5 [19200/162770] Loss: 0.0041 | Recon: 0.0018 | VQ: 0.0023
Epoch 5 [25600/162770] Loss: 0.0049 | Recon: 0.0024 | VQ: 0.0025
Epoch 5 [32000/162770] Loss: 0.0044 | Recon: 0.0020 | VQ: 0.0024
Epoch 5 [38400/162770] Loss: 0.0042 | Recon: 0.0018 | VQ: 0.0024
Epoch 5 [44800/162770] Loss: 0.0040 | Recon: 0.0017 | VQ: 0.0023
Epoch 5 [51200/162770] Loss: 0.0041 | Recon: 0.0019 | VQ: 0.0022
Epoch 5 [57600/162770] Loss: 0.0046 | Recon: 0.0021 | VQ: 0.0025
Epoch 5 [64000/162770] Loss: 0.0041 | Recon: 0.0018 | VQ: 0.0024
Epoch 5 [70400/162770] Loss: 0.0042 | Recon: 0.0018 | VQ: 0.0024
Epoch 5 [76800/162770] Loss: 0.0039 | Recon: 0.0016 | VQ: 0.0023
Epoch 5 [83200/162770] Loss: 0.0042 | Recon: 0.0017 | VQ: 0.0025
Epoch 5 [89600/162770] Loss: 0.0039 | Recon: 0.0018 | VQ: 0.0021
Epoch 5 [96000/162770] Loss: 0.0041 | Recon: 0.0018 | VQ: 0.0023
Epoch 5 [102400/162770] Loss: 0.0039 | Recon: 0.0017 | VQ: 0.0022
Epoch 5 [108800/162770] Loss: 0.0039 | Recon: 0.0016 | VQ: 0.0022
Epoch 5 [115200/162770] Loss: 0.0039 | Recon: 0.0017 | VQ: 0.0022
Epoch 5 [121600/162770] Loss: 0.0040 | Recon: 0.0018 | VQ: 0.0022
Epoch 5 [128000/162770] Loss: 0.0036 | Recon: 0.0016 | VQ: 0.0021
Epoch 5 [134400/162770] Loss: 0.0039 | Recon: 0.0017 | VQ: 0.0021
Epoch 5 [140800/162770] Loss: 0.0039 | Recon: 0.0017 | VQ: 0.0022
Epoch 5 [147200/162770] Loss: 0.0038 | Recon: 0.0017 | VQ: 0.0021
Epoch 5 [153600/162770] Loss: 0.0040 | Recon: 0.0017 | VQ: 0.0023
Epoch 5 [160000/162770] Loss: 0.0039 | Recon: 0.0018 | VQ: 0.0022
====> Epoch: 5 Average Loss: 0.0042 (Recon: 0.0018, VQ: 0.0024)

```

Reasoning: I will visualize the VQ-VAE reconstructions by passing a batch of test images through the trained model and plotting them alongside the original images to qualitatively assess performance.

```

model.eval()

print("Visualizing VQ-VAE Reconstructions...")
with torch.no_grad():
    # Get a batch of data
    data_iter = iter(dataloader)
    images, _ = next(data_iter)
    images = images.to(device)

    # Reconstruct
    _, recon_images, _ = vqvae_model(images)

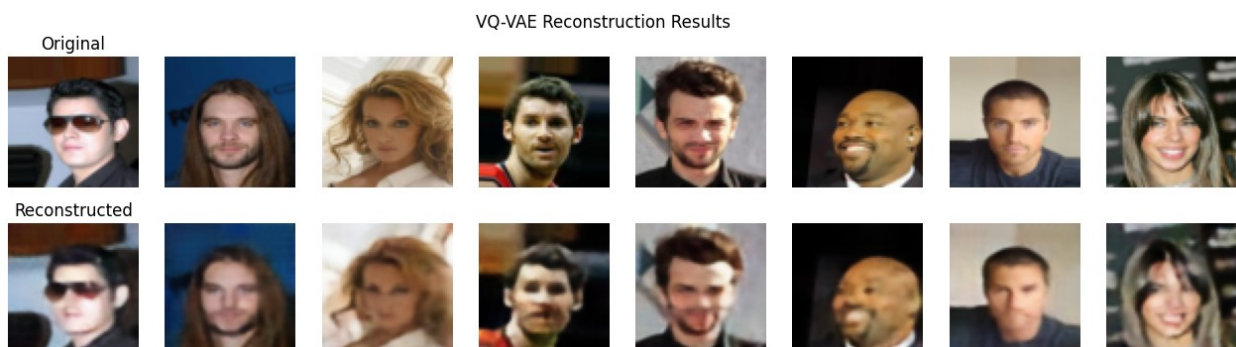
    # Plot
    n = 8
    plt.figure(figsize=(16, 4))
    for i in range(n):
        # Original
        ax = plt.subplot(2, n, i + 1)
        img = images[i].cpu().permute(1, 2, 0).numpy()
        plt.imshow(img)
        plt.axis("off")
        if i == 0: ax.set_title("Original")

        # Reconstructed
        ax = plt.subplot(2, n, i + 1 + n)
        recon_img = recon_images[i].cpu().permute(1, 2, 0).numpy()
        plt.imshow(recon_img)
        plt.axis("off")
        if i == 0: ax.set_title("Reconstructed")

    plt.suptitle("VQ-VAE Reconstruction Results")
    plt.show()

```

Visualizing VQ-VAE Reconstructions...



Part C: PixelCNN Prior

Subtask:

Implement and train a PixelCNN model to learn the prior distribution of the discrete latent codes generated by the VQ-VAE, and use it to sample new images.

Reasoning: Defining the `MaskedConv2d` and `PixelCNN` classes. The `MaskedConv2d` implements the masking logic required for autoregressive generation. The `PixelCNN` uses these layers to model the distribution of latent codes produced by the VQ-VAE.

```
class MaskedConv2d(nn.Conv2d):
    def __init__(self, mask_type, *args, **kwargs):
        super(MaskedConv2d, self).__init__(*args, **kwargs)
        assert mask_type in {'A', 'B'}
        self.register_buffer('mask', self.weight.data.clone())
        _, _, kH, kW = self.weight.size()
        self.mask.fill_(1)
        # Mask logic:
        # Type A: Mask center and future. Used for first layer.
        # Type B: Mask future only (center allowed). Used for
        # subsequent layers.
        self.mask[:, :, kH // 2, kW // 2 + (mask_type == 'B'):] = 0
        self.mask[:, :, kH // 2 + 1:] = 0

    def forward(self, x):
        self.weight.data *= self.mask
        return super(MaskedConv2d, self).forward(x)

class PixelCNN(nn.Module):
    def __init__(self, num_embeddings=512, hidden_dim=64,
num_layers=7):
        super(PixelCNN, self).__init__()
        self.embedding = nn.Embedding(num_embeddings, hidden_dim)

        layers = []
        # First layer: Type A (cannot see current pixel)
        layers.append(MaskedConv2d('A', hidden_dim, hidden_dim,
kernel_size=7, padding=3, bias=False))
        layers.append(nn.BatchNorm2d(hidden_dim))
        layers.append(nn.ReLU(True))

        # Subsequent layers: Type B (can see current pixel from
        # previous layer features)
        for _ in range(num_layers):
            layers.append(MaskedConv2d('B', hidden_dim, hidden_dim,
kernel_size=7, padding=3, bias=False))
            layers.append(nn.BatchNorm2d(hidden_dim))
            layers.append(nn.ReLU(True))
```



```

        self.net = nn.Sequential(*layers)

        # Output layer mapping to codebook size
        self.out = nn.Conv2d(hidden_dim, num_embeddings,
kernel_size=1)

    def forward(self, x):
        # x: (B, H, W) indices
        x = self.embedding(x) # (B, H, W, C)
        x = x.permute(0, 3, 1, 2).contiguous() # (B, C, H, W)
        x = self.net(x)
        return self.out(x)

# Initialize PixelCNN
pixelcnn = PixelCNN(num_embeddings=512, hidden_dim=64).to(device)
print(pixelcnn)

PixelCNN(
  (embedding): Embedding(512, 64)
  (net): Sequential(
    (0): MaskedConv2d(64, 64, kernel_size=(7, 7), stride=(1, 1),
padding=(3, 3), bias=False)
    (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (2): ReLU(inplace=True)
    (3): MaskedConv2d(64, 64, kernel_size=(7, 7), stride=(1, 1),
padding=(3, 3), bias=False)
    (4): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (5): ReLU(inplace=True)
    (6): MaskedConv2d(64, 64, kernel_size=(7, 7), stride=(1, 1),
padding=(3, 3), bias=False)
    (7): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (8): ReLU(inplace=True)
    (9): MaskedConv2d(64, 64, kernel_size=(7, 7), stride=(1, 1),
padding=(3, 3), bias=False)
    (10): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (11): ReLU(inplace=True)
    (12): MaskedConv2d(64, 64, kernel_size=(7, 7), stride=(1, 1),
padding=(3, 3), bias=False)
    (13): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (14): ReLU(inplace=True)
    (15): MaskedConv2d(64, 64, kernel_size=(7, 7), stride=(1, 1),
padding=(3, 3), bias=False)
    (16): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (17): ReLU(inplace=True)

```

```

    (18): MaskedConv2d(64, 64, kernel_size=(7, 7), stride=(1, 1),
padding=(3, 3), bias=False)
    (19): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (20): ReLU(inplace=True)
    (21): MaskedConv2d(64, 64, kernel_size=(7, 7), stride=(1, 1),
padding=(3, 3), bias=False)
    (22): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (23): ReLU(inplace=True)
)
(out): Conv2d(64, 512, kernel_size=(1, 1), stride=(1, 1))
)

```

Reasoning: Implementing the training loop for the PixelCNN model. This involves freezing the VQ-VAE parameters, extracting latent codes from the VQ-VAE, and training the PixelCNN to predict these codes using CrossEntropyLoss.

```

import torch.nn as nn
import torch.optim as optim

# Optimizer for PixelCNN
pixelcnn_optimizer = optim.Adam(pixelcnn.parameters(), lr=1e-3)
criterion = nn.CrossEntropyLoss()

# Ensure VQ-VAE is frozen
for param in vqvae_model.parameters():
    param.requires_grad = False

print("Starting PixelCNN training...")
epochs = 5

for epoch in range(epochs):
    pixelcnn.train()
    train_loss = 0

    for batch_idx, (data, _) in enumerate(dataloader):
        data = data.to(device)

        # 1. Get latent codes from VQ-VAE (frozen)
        with torch.no_grad():
            z = vqvae_model._encoder(data)
            z = vqvae_model._pre_vq_conv(z)
            _, _, _, encoding_indices = vqvae_model._vq_vae(z)

            # Indices shape: (B, H, W, 1) -> (B, H, W)
            target_indices = encoding_indices.view(data.size(0), 16,
16) # 64x64 input -> 16x16 latent

        # 2. Train PixelCNN

```

```

pixelcnn_optimizer.zero_grad()

# Output shape: (B, num_embeddings, H, W)
output = pixelcnn(target_indices)

loss = criterion(output, target_indices)
loss.backward()
pixelcnn_optimizer.step()

train_loss += loss.item()

if batch_idx % 100 == 0:
    print(f'Epoch {epoch+1}
[batch_idx*len(data)]/[len(dataloader.dataset))] Loss:
{loss.item():.4f}')

    print(f'====> Epoch: {epoch+1} Average Loss: {train_loss /
len(dataloader):.4f}')
```

Starting PixelCNN training...

Epoch 1 [0/162770] Loss: 6.3056

Epoch 1 [6400/162770] Loss: 5.3619

Epoch 1 [12800/162770] Loss: 4.9843

Epoch 1 [19200/162770] Loss: 4.6690

Epoch 1 [25600/162770] Loss: 4.3495

Epoch 1 [32000/162770] Loss: 4.1763

Epoch 1 [38400/162770] Loss: 4.1400

Epoch 1 [44800/162770] Loss: 4.1419

Epoch 1 [51200/162770] Loss: 3.9719

Epoch 1 [57600/162770] Loss: 4.0782

Epoch 1 [64000/162770] Loss: 3.9840

Epoch 1 [70400/162770] Loss: 4.0795

Epoch 1 [76800/162770] Loss: 3.8090

Epoch 1 [83200/162770] Loss: 4.0719

Epoch 1 [89600/162770] Loss: 3.8961

Epoch 1 [96000/162770] Loss: 3.8571

Epoch 1 [102400/162770] Loss: 3.8984

Epoch 1 [108800/162770] Loss: 3.9387

Epoch 1 [115200/162770] Loss: 3.7160

Epoch 1 [121600/162770] Loss: 3.7965

Epoch 1 [128000/162770] Loss: 3.7341

Epoch 1 [134400/162770] Loss: 3.8001

Epoch 1 [140800/162770] Loss: 3.6698

Epoch 1 [147200/162770] Loss: 3.8049

Epoch 1 [153600/162770] Loss: 3.7241

Epoch 1 [160000/162770] Loss: 3.7806

====> Epoch: 1 Average Loss: 4.0945

Epoch 2 [0/162770] Loss: 3.8144

Epoch 2 [6400/162770] Loss: 3.6816

Epoch 2 [12800/162770] Loss: 3.7693

```
Epoch 2 [19200/162770] Loss: 3.7383
Epoch 2 [25600/162770] Loss: 3.6796
Epoch 2 [32000/162770] Loss: 3.7679
Epoch 2 [38400/162770] Loss: 3.7868
Epoch 2 [44800/162770] Loss: 3.6917
Epoch 2 [51200/162770] Loss: 3.6391
Epoch 2 [57600/162770] Loss: 3.6901
Epoch 2 [64000/162770] Loss: 3.6888
Epoch 2 [70400/162770] Loss: 3.5599
Epoch 2 [76800/162770] Loss: 3.5761
Epoch 2 [83200/162770] Loss: 3.5913
Epoch 2 [89600/162770] Loss: 3.6295
Epoch 2 [96000/162770] Loss: 3.6688
Epoch 2 [102400/162770] Loss: 3.6378
Epoch 2 [108800/162770] Loss: 3.6840
Epoch 2 [115200/162770] Loss: 3.5695
Epoch 2 [121600/162770] Loss: 3.6691
Epoch 2 [128000/162770] Loss: 3.6492
Epoch 2 [134400/162770] Loss: 3.6058
Epoch 2 [140800/162770] Loss: 3.6496
Epoch 2 [147200/162770] Loss: 3.7291
Epoch 2 [153600/162770] Loss: 3.6688
Epoch 2 [160000/162770] Loss: 3.6578
====> Epoch: 2 Average Loss: 3.6805
Epoch 3 [0/162770] Loss: 3.6676
Epoch 3 [6400/162770] Loss: 3.6186
Epoch 3 [12800/162770] Loss: 3.5774
Epoch 3 [19200/162770] Loss: 3.7015
Epoch 3 [25600/162770] Loss: 3.6601
Epoch 3 [32000/162770] Loss: 3.6192
Epoch 3 [38400/162770] Loss: 3.6528
Epoch 3 [44800/162770] Loss: 3.6090
Epoch 3 [51200/162770] Loss: 3.5579
Epoch 3 [57600/162770] Loss: 3.6444
Epoch 3 [64000/162770] Loss: 3.5780
Epoch 3 [70400/162770] Loss: 3.5743
Epoch 3 [76800/162770] Loss: 3.5243
Epoch 3 [83200/162770] Loss: 3.6746
Epoch 3 [89600/162770] Loss: 3.4816
Epoch 3 [96000/162770] Loss: 3.6349
Epoch 3 [102400/162770] Loss: 3.6434
Epoch 3 [108800/162770] Loss: 3.5971
Epoch 3 [115200/162770] Loss: 3.6082
Epoch 3 [121600/162770] Loss: 3.5361
Epoch 3 [128000/162770] Loss: 3.6219
Epoch 3 [134400/162770] Loss: 3.7090
Epoch 3 [140800/162770] Loss: 3.6221
Epoch 3 [147200/162770] Loss: 3.5857
Epoch 3 [153600/162770] Loss: 3.6024
```

```
Epoch 3 [160000/162770] Loss: 3.5341
====> Epoch: 3 Average Loss: 3.6019
Epoch 4 [0/162770] Loss: 3.7807
Epoch 4 [6400/162770] Loss: 3.5978
Epoch 4 [12800/162770] Loss: 3.5125
Epoch 4 [19200/162770] Loss: 3.6372
Epoch 4 [25600/162770] Loss: 3.5296
Epoch 4 [32000/162770] Loss: 3.5212
Epoch 4 [38400/162770] Loss: 3.6259
Epoch 4 [44800/162770] Loss: 3.5414
Epoch 4 [51200/162770] Loss: 3.5458
Epoch 4 [57600/162770] Loss: 3.5239
Epoch 4 [64000/162770] Loss: 3.5383
Epoch 4 [70400/162770] Loss: 3.6888
Epoch 4 [76800/162770] Loss: 3.5274
Epoch 4 [83200/162770] Loss: 3.5495
Epoch 4 [89600/162770] Loss: 3.5132
Epoch 4 [96000/162770] Loss: 3.6720
Epoch 4 [102400/162770] Loss: 3.5866
Epoch 4 [108800/162770] Loss: 3.6072
Epoch 4 [115200/162770] Loss: 3.5478
Epoch 4 [121600/162770] Loss: 3.5417
Epoch 4 [128000/162770] Loss: 3.5603
Epoch 4 [134400/162770] Loss: 3.6386
Epoch 4 [140800/162770] Loss: 3.5974
Epoch 4 [147200/162770] Loss: 3.5399
Epoch 4 [153600/162770] Loss: 3.5392
Epoch 4 [160000/162770] Loss: 3.4644
====> Epoch: 4 Average Loss: 3.5658
Epoch 5 [0/162770] Loss: 3.5762
Epoch 5 [6400/162770] Loss: 3.5215
Epoch 5 [12800/162770] Loss: 3.4048
Epoch 5 [19200/162770] Loss: 3.5348
Epoch 5 [25600/162770] Loss: 3.4545
Epoch 5 [32000/162770] Loss: 3.5455
Epoch 5 [38400/162770] Loss: 3.5166
Epoch 5 [44800/162770] Loss: 3.4410
Epoch 5 [51200/162770] Loss: 3.3902
Epoch 5 [57600/162770] Loss: 3.5962
Epoch 5 [64000/162770] Loss: 3.5587
Epoch 5 [70400/162770] Loss: 3.5281
Epoch 5 [76800/162770] Loss: 3.5780
Epoch 5 [83200/162770] Loss: 3.4882
Epoch 5 [89600/162770] Loss: 3.5981
Epoch 5 [96000/162770] Loss: 3.5308
Epoch 5 [102400/162770] Loss: 3.5020
Epoch 5 [108800/162770] Loss: 3.5338
Epoch 5 [115200/162770] Loss: 3.5499
Epoch 5 [121600/162770] Loss: 3.5495
```

```
Epoch 5 [128000/162770] Loss: 3.5870
Epoch 5 [134400/162770] Loss: 3.6450
Epoch 5 [140800/162770] Loss: 3.5525
Epoch 5 [147200/162770] Loss: 3.4278
Epoch 5 [153600/162770] Loss: 3.4484
Epoch 5 [160000/162770] Loss: 3.4715
====> Epoch: 5 Average Loss: 3.5412
```

Reasoning: Implementing the sampling logic for the PixelCNN. This involves iteratively generating latent codes pixel-by-pixel, using the trained PixelCNN to predict the distribution of the next token based on previous ones. After generating the full latent map, I will use the VQ-VAE decoder to synthesize the corresponding images and visualize the results.

```
def generate_samples(num=16):
    pixelcnn.eval()
    # Initialize latent array (batch, H, W)
    # VQ-VAE downsamples 64x64 by factor 4 -> 16x16 latents
    samples = torch.zeros((num, 16, 16), dtype=torch.long).to(device)

    print("Sampling Latent Codes (Autoregressive)...")
    with torch.no_grad():
        for i in range(16):
            for j in range(16):
                output = pixelcnn(samples)
                # Get logits for the specific pixel (i, j)
                # output shape: (B, num_embeddings, H, W)
                logits = output[:, :, i, j]
                probs = F.softmax(logits, dim=1)

                # Sample from the distribution
                samples[:, i, j] = torch.multinomial(probs,
1).squeeze(1)

    # Decode samples using VQ-VAE
    print("Decoding samples...")
    with torch.no_grad():
        # Convert indices to vectors
        # samples: (B, 16, 16)
        # embedding weight: (num_embeddings, embedding_dim)

        # Flatten samples to index embedding
        flat_samples = samples.view(-1)

        # Get embeddings: (B*H*W, C)
        quantized_flat = vqvae_model._vq_vae._embedding(flat_samples)

        # Reshape to (B, C, H, W) for decoder
        # Note: Embedding output is (B*H*W, C), need (B, C, H, W)
        quantized = quantized_flat.view(num, 16, 16, 64).permute(0, 3,
```

```
1, 2).contiguous()

    generated_images = vqvae_model._decoder(quantized)

    # Visualize
    plt.figure(figsize=(8, 8))
    grid_img = torchvision.utils.make_grid(generated_images.cpu(),
nrow=4, padding=2, normalize=False)
    plt.imshow(grid_img.permute(1, 2, 0).numpy())
    plt.axis("off")
    plt.title("PixelCNN + VQ-VAE Generated Samples")
    plt.show()

# Execute generation
generate_samples(num=16)

Sampling Latent Codes (Autoregressive)...
Decoding samples...
```

PixelCNN + VQ-VAE Generated Samples



Part D: GAN Implementation and Training

Subtask:

Explain the GAN Minimax objective, implement the Generator and Discriminator architectures, and train the GAN on the CelebA dataset.

1. GAN Minimax Objective

Generative Adversarial Networks (GANs) consist of two neural networks, a **Generator (G)** and a **Discriminator (D)**, which compete in a zero-sum game.

- **Generator (G):** Tries to capture the data distribution p_{data} by mapping a latent noise vector z (sampled from a prior p_z , usually Gaussian) to the data space ($G(z)$). Its goal is to generate images indistinguishable from real data.
- **Discriminator (D):** Tries to distinguish between real data samples x and fake samples $G(z)$ produced by the generator. It outputs a probability value indicating the likelihood that the input is real.

The training process is formulated as a minimax game with the following objective function:

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)} [\log D(x)] + E_{z \sim p_z(z)} [\log (1 - D(G(z)))]$$

- The **Discriminator** maximizes this objective: it wants $D(x)$ close to 1 (log term $\rightarrow 0$) and $D(G(z))$ close to 0 (log term $\rightarrow 0$).
- The **Generator** minimizes this objective: it wants $D(G(z))$ close to 1 (log term $\rightarrow -\infty$).

Reasoning: Implementing the Generator and Discriminator classes suitable for a DCGAN architecture on 64x64 images, along with the standard weight initialization function.

```
import torch.nn as nn

# Custom weights initialization called on netG and netD
def weights_init(m):
    classname = m.__class__.__name__
    if classname.find('Conv') != -1:
        nn.init.normal_(m.weight.data, 0.0, 0.02)
    elif classname.find('BatchNorm') != -1:
        nn.init.normal_(m.weight.data, 1.0, 0.02)
        nn.init.constant_(m.bias.data, 0)

# Generator Code
class Generator(nn.Module):
    def __init__(self, nz=100, ngf=64, nc=3):
        super(Generator, self).__init__()
        self.main = nn.Sequential(
            # input is Z, going into a convolution
            nn.ConvTranspose2d(nz, ngf * 8, 4, 1, 0, bias=False),
            nn.BatchNorm2d(ngf * 8),
            nn.ReLU(True),
            # state size. (ngf*8) x 4 x 4
            nn.ConvTranspose2d(ngf * 8, ngf * 4, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf * 4),
            nn.ReLU(True),
            # state size. (ngf*4) x 8 x 8
            nn.ConvTranspose2d(ngf * 4, ngf * 2, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ngf * 2),
            nn.ReLU(True),
            # state size. (ngf*2) x 16 x 16
            nn.ConvTranspose2d(ngf * 2, ngf, 4, 2, 1, bias=False),
```

```

        nn.BatchNorm2d(ngf),
        nn.ReLU(True),
        # state size. (ngf) x 32 x 32
        nn.ConvTranspose2d(ngf, nc, 4, 2, 1, bias=False),
        nn.Tanh()
        # state size. (nc) x 64 x 64
    )

    def forward(self, input):
        return self.main(input)

# Discriminator Code
class Discriminator(nn.Module):
    def __init__(self, nc=3, ndf=64):
        super(Discriminator, self).__init__()
        self.main = nn.Sequential(
            # input is (nc) x 64 x 64
            nn.Conv2d(nc, ndf, 4, 2, 1, bias=False),
            nn.LeakyReLU(0.2, inplace=True),
            # state size. (ndf) x 32 x 32
            nn.Conv2d(ndf, ndf * 2, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ndf * 2),
            nn.LeakyReLU(0.2, inplace=True),
            # state size. (ndf*2) x 16 x 16
            nn.Conv2d(ndf * 2, ndf * 4, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ndf * 4),
            nn.LeakyReLU(0.2, inplace=True),
            # state size. (ndf*4) x 8 x 8
            nn.Conv2d(ndf * 4, ndf * 8, 4, 2, 1, bias=False),
            nn.BatchNorm2d(ndf * 8),
            nn.LeakyReLU(0.2, inplace=True),
            # state size. (ndf*8) x 4 x 4
            nn.Conv2d(ndf * 8, 1, 4, 1, 0, bias=False),
            nn.Sigmoid()
        )

    def forward(self, input):
        return self.main(input)

# Initialize Models
nz = 100
netG = Generator(nz=nz).to(device)
netD = Discriminator().to(device)

# Apply weights
netG.apply(weights_init)
netD.apply(weights_init)

print(netG)
print(netD)

```

```

Generator(
  (main): Sequential(
    (0): ConvTranspose2d(100, 512, kernel_size=(4, 4), stride=(1, 1),
bias=False)
    (1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (2): ReLU(inplace=True)
    (3): ConvTranspose2d(512, 256, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1), bias=False)
    (4): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (5): ReLU(inplace=True)
    (6): ConvTranspose2d(256, 128, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1), bias=False)
    (7): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (8): ReLU(inplace=True)
    (9): ConvTranspose2d(128, 64, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1), bias=False)
    (10): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (11): ReLU(inplace=True)
    (12): ConvTranspose2d(64, 3, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1), bias=False)
    (13): Tanh()
  )
)
Discriminator(
  (main): Sequential(
    (0): Conv2d(3, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1,
1), bias=False)
    (1): LeakyReLU(negative_slope=0.2, inplace=True)
    (2): Conv2d(64, 128, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1), bias=False)
    (3): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (4): LeakyReLU(negative_slope=0.2, inplace=True)
    (5): Conv2d(128, 256, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1), bias=False)
    (6): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (7): LeakyReLU(negative_slope=0.2, inplace=True)
    (8): Conv2d(256, 512, kernel_size=(4, 4), stride=(2, 2),
padding=(1, 1), bias=False)
    (9): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True,
track_running_stats=True)
    (10): LeakyReLU(negative_slope=0.2, inplace=True)
    (11): Conv2d(512, 1, kernel_size=(4, 4), stride=(1, 1),
bias=False)
    (12): Sigmoid()
  )
)

```

```
)  
)
```

Reasoning: Implementing the GAN training loop. This involves defining the loss function (Binary Cross Entropy), initializing the optimizers for both the Generator and Discriminator, creating fixed noise for consistent visualization, and iterating through the dataset for 5 epochs to update both networks according to the Minimax objective.

```
import torch.optim as optim  
  
# Initialize BCELoss function  
criterion = nn.BCELoss()  
  
# Create batch of latent vectors that we will use to visualize  
# the progression of the generator  
fixed_noise = torch.randn(64, nz, 1, 1, device=device)  
  
# Establish convention for real and fake labels during training  
real_label = 1.  
fake_label = 0.  
  
# Setup Adam optimizers for both G and D  
optimizerD = optim.Adam(netD.parameters(), lr=0.0002, betas=(0.5,  
0.999))  
optimizerG = optim.Adam(netG.parameters(), lr=0.0002, betas=(0.5,  
0.999))  
  
print("Starting GAN Training Loop...")  
# For 5 epochs  
num_epochs = 5  
  
for epoch in range(num_epochs):  
    for i, (data, _) in enumerate(dataloader):  
        #####  
        # (1) Update D network: maximize log(D(x)) + log(1 - D(G(z)))  
        #####  
        ## Train with all-real batch  
        netD.zero_grad()  
        # Format batch  
        real_cpu = data.to(device)  
        b_size = real_cpu.size(0)  
        label = torch.full((b_size,), real_label, dtype=torch.float,  
device=device)  
  
        # Forward pass real batch through D  
        output = netD(real_cpu).view(-1)  
        # Calculate loss on all-real batch  
        errD_real = criterion(output, label)  
        # Calculate gradients for D in backward pass  
        errD_real.backward()
```

```

D_x = output.mean().item()

## Train with all-fake batch
# Generate batch of latent vectors
noise = torch.randn(b_size, nz, 1, 1, device=device)
# Generate fake image batch with G
fake = netG(noise)
label.fill_(fake_label)

# Classify all fake batch with D
output = netD(fake.detach()).view(-1)
# Calculate D's loss on the all-fake batch
errD_fake = criterion(output, label)
# Calculate the gradients for this batch, accumulated (summed)
with previous gradients
errD_fake.backward()
D_G_z1 = output.mean().item()

# Compute error of D as sum over the fake and the real batches
errD = errD_real + errD_fake
# Update D
optimizerD.step()

#####
# (2) Update G network: maximize log(D(G(z)))
#####
netG.zero_grad()
label.fill_(real_label) # fake labels are real for generator

cost

# Since we just updated D, perform another forward pass of
all-fake batch through D
output = netD(fake).view(-1)
# Calculate G's loss based on this output
errG = criterion(output, label)
# Calculate gradients for G
errG.backward()
D_G_z2 = output.mean().item()
# Update G
optimizerG.step()

# Output training stats
if i % 100 == 0:
    print(f'[{epoch+1}/{num_epochs}]
[{i*len(data)}/{len(dataloader.dataset)}] '
          f'Loss_D: {errD.item():.4f} Loss_G:
{errG.item():.4f} '
          f'D(x): {D_x:.4f} D(G(z)): {D_G_z1:.4f} /
{D_G_z2:.4f}')

```

```
print(f'====> Epoch: {epoch+1} Complete')
```

Starting GAN Training Loop...

```
[1/5][0/162770] Loss_D: 1.7677 Loss_G: 2.3404 D(x): 0.3362 D(G(z)): 0.3574 / 0.1152
```

```
[1/5][6400/162770] Loss_D: 0.4018 Loss_G: 14.6204 D(x): 0.7867 D(G(z)): 0.0000 / 0.0000
```

```
[1/5][12800/162770] Loss_D: 2.1843 Loss_G: 3.2225 D(x): 0.2392 D(G(z)): 0.0069 / 0.0800
```

```
[1/5][19200/162770] Loss_D: 1.5655 Loss_G: 1.7854 D(x): 0.3123 D(G(z)): 0.0153 / 0.2547
```

```
[1/5][25600/162770] Loss_D: 0.7429 Loss_G: 3.2662 D(x): 0.6661 D(G(z)): 0.1717 / 0.0584
```

```
[1/5][32000/162770] Loss_D: 0.9726 Loss_G: 2.1906 D(x): 0.5249 D(G(z)): 0.0290 / 0.1587
```

```
[1/5][38400/162770] Loss_D: 0.6071 Loss_G: 2.1159 D(x): 0.7119 D(G(z)): 0.1442 / 0.1751
```

```
[1/5][44800/162770] Loss_D: 0.6347 Loss_G: 3.1025 D(x): 0.6669 D(G(z)): 0.0510 / 0.0701
```

```
[1/5][51200/162770] Loss_D: 0.9790 Loss_G: 3.1251 D(x): 0.5041 D(G(z)): 0.0079 / 0.0780
```

```
[1/5][57600/162770] Loss_D: 1.3287 Loss_G: 1.6301 D(x): 0.4560 D(G(z)): 0.1356 / 0.2791
```

```
[1/5][64000/162770] Loss_D: 0.6616 Loss_G: 2.9858 D(x): 0.7044 D(G(z)): 0.2048 / 0.0704
```

```
[1/5][70400/162770] Loss_D: 0.6433 Loss_G: 2.6979 D(x): 0.6504 D(G(z)): 0.0639 / 0.1135
```

```
[1/5][76800/162770] Loss_D: 0.9467 Loss_G: 5.6566 D(x): 0.8696 D(G(z)): 0.4851 / 0.0080
```

```
[1/5][83200/162770] Loss_D: 0.8255 Loss_G: 3.5753 D(x): 0.5658 D(G(z)): 0.0394 / 0.0818
```

```
[1/5][89600/162770] Loss_D: 0.3997 Loss_G: 3.8164 D(x): 0.8697 D(G(z)): 0.1999 / 0.0378
```

```
[1/5][96000/162770] Loss_D: 0.2826 Loss_G: 3.8212 D(x): 0.8702 D(G(z)): 0.1135 / 0.0324
```

```
[1/5][102400/162770] Loss_D: 0.6728 Loss_G: 4.9153 D(x): 0.8756 D(G(z)): 0.3622 / 0.0114
```

```
[1/5][108800/162770] Loss_D: 0.5475 Loss_G: 3.6113 D(x): 0.8784 D(G(z)): 0.2813 / 0.0384
```

```
[1/5][115200/162770] Loss_D: 0.3903 Loss_G: 3.0152 D(x): 0.8416 D(G(z)): 0.1628 / 0.0751
```

```
[1/5][121600/162770] Loss_D: 0.6352 Loss_G: 2.5020 D(x): 0.6429 D(G(z)): 0.0683 / 0.1191
```

```
[1/5][128000/162770] Loss_D: 1.1389 Loss_G: 2.3001 D(x): 0.4568 D(G(z)): 0.0123 / 0.1553
```

```
[1/5][134400/162770] Loss_D: 0.6106 Loss_G: 2.6575 D(x): 0.7319 D(G(z)): 0.1952 / 0.1152
```

```
[1/5][140800/162770] Loss_D: 1.3755 Loss_G: 0.9806 D(x): 0.3307 D(G(z)): 0.0262 / 0.4466
```

[1/5][147200/162770] Loss_D: 0.6773 Loss_G: 2.0237 D(x): 0.6667
D(G(z)): 0.1219 / 0.1837
[1/5][153600/162770] Loss_D: 0.6312 Loss_G: 2.2395 D(x): 0.6687
D(G(z)): 0.0956 / 0.1605
[1/5][160000/162770] Loss_D: 0.9334 Loss_G: 4.2466 D(x): 0.8408
D(G(z)): 0.4689 / 0.0209
====> Epoch: 1 Complete
[2/5][0/162770] Loss_D: 0.9540 Loss_G: 1.5424 D(x): 0.5047 D(G(z)): 0.0482 / 0.2839
[2/5][6400/162770] Loss_D: 0.6911 Loss_G: 2.7097 D(x): 0.7704 D(G(z)): 0.2568 / 0.1002
[2/5][12800/162770] Loss_D: 0.3837 Loss_G: 3.5113 D(x): 0.8111
D(G(z)): 0.1200 / 0.0499
[2/5][19200/162770] Loss_D: 0.7038 Loss_G: 3.1060 D(x): 0.7762
D(G(z)): 0.2974 / 0.0621
[2/5][25600/162770] Loss_D: 0.8307 Loss_G: 1.6493 D(x): 0.5331
D(G(z)): 0.0632 / 0.2529
[2/5][32000/162770] Loss_D: 0.5941 Loss_G: 2.2034 D(x): 0.7753
D(G(z)): 0.2330 / 0.1404
[2/5][38400/162770] Loss_D: 0.3812 Loss_G: 3.6494 D(x): 0.8343
D(G(z)): 0.1492 / 0.0395
[2/5][44800/162770] Loss_D: 1.7806 Loss_G: 6.6250 D(x): 0.9744
D(G(z)): 0.7475 / 0.0049
[2/5][51200/162770] Loss_D: 0.6751 Loss_G: 1.2655 D(x): 0.6401
D(G(z)): 0.1184 / 0.3147
[2/5][57600/162770] Loss_D: 0.5420 Loss_G: 3.3154 D(x): 0.8005
D(G(z)): 0.2263 / 0.0495
[2/5][64000/162770] Loss_D: 0.5851 Loss_G: 2.3980 D(x): 0.6648
D(G(z)): 0.0812 / 0.1233
[2/5][70400/162770] Loss_D: 0.5990 Loss_G: 2.4171 D(x): 0.6607
D(G(z)): 0.0897 / 0.1339
[2/5][76800/162770] Loss_D: 0.5501 Loss_G: 2.6679 D(x): 0.8836
D(G(z)): 0.3142 / 0.0896
[2/5][83200/162770] Loss_D: 0.6607 Loss_G: 4.0661 D(x): 0.9220
D(G(z)): 0.3996 / 0.0248
[2/5][89600/162770] Loss_D: 0.4774 Loss_G: 2.9337 D(x): 0.7942
D(G(z)): 0.1835 / 0.0671
[2/5][96000/162770] Loss_D: 0.5054 Loss_G: 1.6566 D(x): 0.7461
D(G(z)): 0.1524 / 0.2200
[2/5][102400/162770] Loss_D: 0.7605 Loss_G: 3.1298 D(x): 0.6918
D(G(z)): 0.2680 / 0.0625
[2/5][108800/162770] Loss_D: 1.1262 Loss_G: 4.3799 D(x): 0.9326
D(G(z)): 0.5773 / 0.0196
[2/5][115200/162770] Loss_D: 0.5469 Loss_G: 3.0833 D(x): 0.8238
D(G(z)): 0.2547 / 0.0641
[2/5][121600/162770] Loss_D: 1.3867 Loss_G: 0.5124 D(x): 0.3234
D(G(z)): 0.0198 / 0.6477
[2/5][128000/162770] Loss_D: 0.7833 Loss_G: 2.6250 D(x): 0.7890
D(G(z)): 0.3706 / 0.0899

[2/5][134400/162770] Loss_D: 1.0150 Loss_G: 5.1455 D(x): 0.8556
D(G(z)): 0.5247 / 0.0082
[2/5][140800/162770] Loss_D: 0.4322 Loss_G: 3.4651 D(x): 0.8751
D(G(z)): 0.2320 / 0.0410
[2/5][147200/162770] Loss_D: 2.2496 Loss_G: 5.5357 D(x): 0.9724
D(G(z)): 0.8436 / 0.0096
[2/5][153600/162770] Loss_D: 2.0070 Loss_G: 5.0406 D(x): 0.8024
D(G(z)): 0.7470 / 0.0122
[2/5][160000/162770] Loss_D: 0.8870 Loss_G: 2.4375 D(x): 0.6527
D(G(z)): 0.2911 / 0.1107
====> Epoch: 2 Complete
[3/5][0/162770] Loss_D: 0.7965 Loss_G: 3.4107 D(x): 0.8023 D(G(z)):
0.3866 / 0.0536
[3/5][6400/162770] Loss_D: 2.6003 Loss_G: 0.8868 D(x): 0.1318 D(G(z)):
0.0253 / 0.5205
[3/5][12800/162770] Loss_D: 1.1392 Loss_G: 4.6985 D(x): 0.9625
D(G(z)): 0.6164 / 0.0129
[3/5][19200/162770] Loss_D: 0.5065 Loss_G: 2.7329 D(x): 0.8458
D(G(z)): 0.2623 / 0.0779
[3/5][25600/162770] Loss_D: 0.7447 Loss_G: 2.1140 D(x): 0.6839
D(G(z)): 0.2314 / 0.1729
[3/5][32000/162770] Loss_D: 0.6620 Loss_G: 2.8437 D(x): 0.7551
D(G(z)): 0.2755 / 0.0844
[3/5][38400/162770] Loss_D: 0.6223 Loss_G: 2.1345 D(x): 0.6659
D(G(z)): 0.1226 / 0.1553
[3/5][44800/162770] Loss_D: 0.8868 Loss_G: 1.1639 D(x): 0.5256
D(G(z)): 0.1292 / 0.3758
[3/5][51200/162770] Loss_D: 0.7097 Loss_G: 1.2146 D(x): 0.6896
D(G(z)): 0.2205 / 0.3446
[3/5][57600/162770] Loss_D: 0.8079 Loss_G: 1.1588 D(x): 0.5567
D(G(z)): 0.1204 / 0.3483
[3/5][64000/162770] Loss_D: 0.7269 Loss_G: 2.4197 D(x): 0.7276
D(G(z)): 0.2841 / 0.1016
[3/5][70400/162770] Loss_D: 0.8612 Loss_G: 1.3977 D(x): 0.5260
D(G(z)): 0.1157 / 0.2966
[3/5][76800/162770] Loss_D: 0.9551 Loss_G: 4.3219 D(x): 0.9478
D(G(z)): 0.5292 / 0.0210
[3/5][83200/162770] Loss_D: 0.6402 Loss_G: 2.0149 D(x): 0.6188
D(G(z)): 0.0719 / 0.1749
[3/5][89600/162770] Loss_D: 0.6264 Loss_G: 2.0937 D(x): 0.7204
D(G(z)): 0.2061 / 0.1667
[3/5][96000/162770] Loss_D: 0.5883 Loss_G: 1.7667 D(x): 0.7174
D(G(z)): 0.1594 / 0.2098
[3/5][102400/162770] Loss_D: 0.7987 Loss_G: 1.9268 D(x): 0.6330
D(G(z)): 0.2057 / 0.1827
[3/5][108800/162770] Loss_D: 0.7369 Loss_G: 2.7450 D(x): 0.7069
D(G(z)): 0.2740 / 0.0843
[3/5][115200/162770] Loss_D: 0.4775 Loss_G: 2.5008 D(x): 0.7898
D(G(z)): 0.1789 / 0.1056

[3/5][121600/162770] Loss_D: 1.0023 Loss_G: 0.9910 D(x): 0.4668
D(G(z)): 0.1010 / 0.4399
[3/5][128000/162770] Loss_D: 0.4830 Loss_G: 2.3650 D(x): 0.7056
D(G(z)): 0.0820 / 0.1214
[3/5][134400/162770] Loss_D: 0.8278 Loss_G: 2.5202 D(x): 0.6865
D(G(z)): 0.2965 / 0.1024
[3/5][140800/162770] Loss_D: 1.3002 Loss_G: 3.4758 D(x): 0.9016
D(G(z)): 0.5728 / 0.0588
[3/5][147200/162770] Loss_D: 0.5058 Loss_G: 2.3804 D(x): 0.7742
D(G(z)): 0.1869 / 0.1141
[3/5][153600/162770] Loss_D: 0.7036 Loss_G: 2.3677 D(x): 0.7480
D(G(z)): 0.2823 / 0.1199
[3/5][160000/162770] Loss_D: 0.5245 Loss_G: 2.3511 D(x): 0.6881
D(G(z)): 0.1018 / 0.1155
====> Epoch: 3 Complete
[4/5][0/162770] Loss_D: 0.6429 Loss_G: 2.5850 D(x): 0.8061 D(G(z)): 0.2867 / 0.1019
[4/5][6400/162770] Loss_D: 1.2472 Loss_G: 4.9615 D(x): 0.9768 D(G(z)): 0.6469 / 0.0116
[4/5][12800/162770] Loss_D: 0.8245 Loss_G: 2.6044 D(x): 0.7249
D(G(z)): 0.3395 / 0.0999
[4/5][19200/162770] Loss_D: 0.4419 Loss_G: 2.4562 D(x): 0.7740
D(G(z)): 0.1397 / 0.1124
[4/5][25600/162770] Loss_D: 1.6483 Loss_G: 7.0062 D(x): 0.9627
D(G(z)): 0.7125 / 0.0021
[4/5][32000/162770] Loss_D: 0.4728 Loss_G: 3.2009 D(x): 0.8723
D(G(z)): 0.2564 / 0.0525
[4/5][38400/162770] Loss_D: 1.6413 Loss_G: 0.1555 D(x): 0.2725
D(G(z)): 0.0375 / 0.8723
[4/5][44800/162770] Loss_D: 0.6451 Loss_G: 4.3793 D(x): 0.9247
D(G(z)): 0.3988 / 0.0182
[4/5][51200/162770] Loss_D: 0.4714 Loss_G: 2.8461 D(x): 0.7746
D(G(z)): 0.1522 / 0.0828
[4/5][57600/162770] Loss_D: 0.8263 Loss_G: 1.4427 D(x): 0.6139
D(G(z)): 0.2176 / 0.2725
[4/5][64000/162770] Loss_D: 0.8599 Loss_G: 1.0578 D(x): 0.5188
D(G(z)): 0.1013 / 0.3966
[4/5][70400/162770] Loss_D: 0.6536 Loss_G: 2.4910 D(x): 0.6915
D(G(z)): 0.1648 / 0.1087
[4/5][76800/162770] Loss_D: 0.4798 Loss_G: 3.3007 D(x): 0.8223
D(G(z)): 0.2121 / 0.0498
[4/5][83200/162770] Loss_D: 0.4108 Loss_G: 2.8827 D(x): 0.8562
D(G(z)): 0.2024 / 0.0720
[4/5][89600/162770] Loss_D: 1.0375 Loss_G: 5.1939 D(x): 0.9457
D(G(z)): 0.5755 / 0.0081
[4/5][96000/162770] Loss_D: 0.6310 Loss_G: 1.9076 D(x): 0.6588
D(G(z)): 0.1549 / 0.1849
[4/5][102400/162770] Loss_D: 0.5142 Loss_G: 2.1098 D(x): 0.6790
D(G(z)): 0.0783 / 0.1512

[4/5][108800/162770] Loss_D: 1.4958 Loss_G: 4.0788 D(x): 0.8254
D(G(z)): 0.6368 / 0.0321
[4/5][115200/162770] Loss_D: 0.9671 Loss_G: 1.4123 D(x): 0.4684
D(G(z)): 0.0776 / 0.2902
[4/5][121600/162770] Loss_D: 0.7052 Loss_G: 3.6756 D(x): 0.8361
D(G(z)): 0.3633 / 0.0333
[4/5][128000/162770] Loss_D: 0.5203 Loss_G: 3.3896 D(x): 0.9085
D(G(z)): 0.3117 / 0.0473
[4/5][134400/162770] Loss_D: 0.6266 Loss_G: 2.4394 D(x): 0.7637
D(G(z)): 0.2405 / 0.1167
[4/5][140800/162770] Loss_D: 1.0483 Loss_G: 5.0221 D(x): 0.9127
D(G(z)): 0.5495 / 0.0133
[4/5][147200/162770] Loss_D: 0.5306 Loss_G: 2.1234 D(x): 0.7526
D(G(z)): 0.1823 / 0.1617
[4/5][153600/162770] Loss_D: 0.7627 Loss_G: 2.2333 D(x): 0.6765
D(G(z)): 0.2341 / 0.1385
[4/5][160000/162770] Loss_D: 0.4239 Loss_G: 3.0566 D(x): 0.7726
D(G(z)): 0.1225 / 0.0654
====> Epoch: 4 Complete
[5/5][0/162770] Loss_D: 2.2185 Loss_G: 5.2375 D(x): 0.9913 D(G(z)): 0.8374 / 0.0149
[5/5][6400/162770] Loss_D: 0.7831 Loss_G: 1.2587 D(x): 0.5320 D(G(z)): 0.0489 / 0.3478
[5/5][12800/162770] Loss_D: 2.0696 Loss_G: 0.4291 D(x): 0.2159
D(G(z)): 0.0497 / 0.7108
[5/5][19200/162770] Loss_D: 0.5250 Loss_G: 3.1823 D(x): 0.8304
D(G(z)): 0.2431 / 0.0559
[5/5][25600/162770] Loss_D: 0.6174 Loss_G: 3.4366 D(x): 0.9190
D(G(z)): 0.3732 / 0.0438
[5/5][32000/162770] Loss_D: 0.6021 Loss_G: 3.6450 D(x): 0.9023
D(G(z)): 0.3523 / 0.0366
[5/5][38400/162770] Loss_D: 0.9416 Loss_G: 4.4928 D(x): 0.8810
D(G(z)): 0.4585 / 0.0202
[5/5][44800/162770] Loss_D: 0.3680 Loss_G: 2.9566 D(x): 0.8580
D(G(z)): 0.1733 / 0.0673
[5/5][51200/162770] Loss_D: 0.4339 Loss_G: 2.2891 D(x): 0.8237
D(G(z)): 0.1915 / 0.1281
[5/5][57600/162770] Loss_D: 0.7531 Loss_G: 2.1199 D(x): 0.6398
D(G(z)): 0.1507 / 0.1631
[5/5][64000/162770] Loss_D: 0.4478 Loss_G: 2.3071 D(x): 0.8009
D(G(z)): 0.1797 / 0.1232
[5/5][70400/162770] Loss_D: 0.6514 Loss_G: 3.3015 D(x): 0.8354
D(G(z)): 0.3223 / 0.0546
[5/5][76800/162770] Loss_D: 0.6130 Loss_G: 1.4745 D(x): 0.6252
D(G(z)): 0.0648 / 0.2909
[5/5][83200/162770] Loss_D: 0.7988 Loss_G: 2.4548 D(x): 0.7506
D(G(z)): 0.3258 / 0.1118
[5/5][89600/162770] Loss_D: 0.4738 Loss_G: 2.2355 D(x): 0.7378
D(G(z)): 0.1126 / 0.1486

```

[5/5][96000/162770] Loss_D: 0.5282 Loss_G: 3.1784 D(x): 0.8391
D(G(z)): 0.2535 / 0.0528
[5/5][102400/162770] Loss_D: 0.6909 Loss_G: 4.8949 D(x): 0.9733
D(G(z)): 0.4366 / 0.0110
[5/5][108800/162770] Loss_D: 0.4816 Loss_G: 3.7409 D(x): 0.8336
D(G(z)): 0.2223 / 0.0309
[5/5][115200/162770] Loss_D: 0.4159 Loss_G: 2.3414 D(x): 0.7967
D(G(z)): 0.1475 / 0.1300
[5/5][121600/162770] Loss_D: 0.5645 Loss_G: 2.8626 D(x): 0.8339
D(G(z)): 0.2904 / 0.0667
[5/5][128000/162770] Loss_D: 0.3984 Loss_G: 3.1667 D(x): 0.8030
D(G(z)): 0.1387 / 0.0628
[5/5][134400/162770] Loss_D: 0.6217 Loss_G: 4.0600 D(x): 0.8796
D(G(z)): 0.3587 / 0.0211
[5/5][140800/162770] Loss_D: 0.3786 Loss_G: 2.1074 D(x): 0.8027
D(G(z)): 0.1160 / 0.1526
[5/5][147200/162770] Loss_D: 0.7217 Loss_G: 1.0774 D(x): 0.5673
D(G(z)): 0.0453 / 0.4054
[5/5][153600/162770] Loss_D: 0.3907 Loss_G: 2.8084 D(x): 0.8097
D(G(z)): 0.1242 / 0.0882
[5/5][160000/162770] Loss_D: 0.5151 Loss_G: 2.7812 D(x): 0.8451
D(G(z)): 0.2659 / 0.0802
====> Epoch: 5 Complete

```

Reasoning: Visualizing the results of the GAN training by using the generator `netG` to produce images from the `fixed_noise` vector. This allows us to qualitatively assess the realism of the generated faces.

```

import torchvision.utils as vutils

# Visualize the results
print("Visualizing GAN Generated Images...")
with torch.no_grad():
    # Generate fake images using the fixed noise vector
    fake_images = netG(fixed_noise).detach().cpu()

    # Plot
    plt.figure(figsize=(8, 8))
    plt.axis("off")
    plt.title("Generated Images (GAN)")

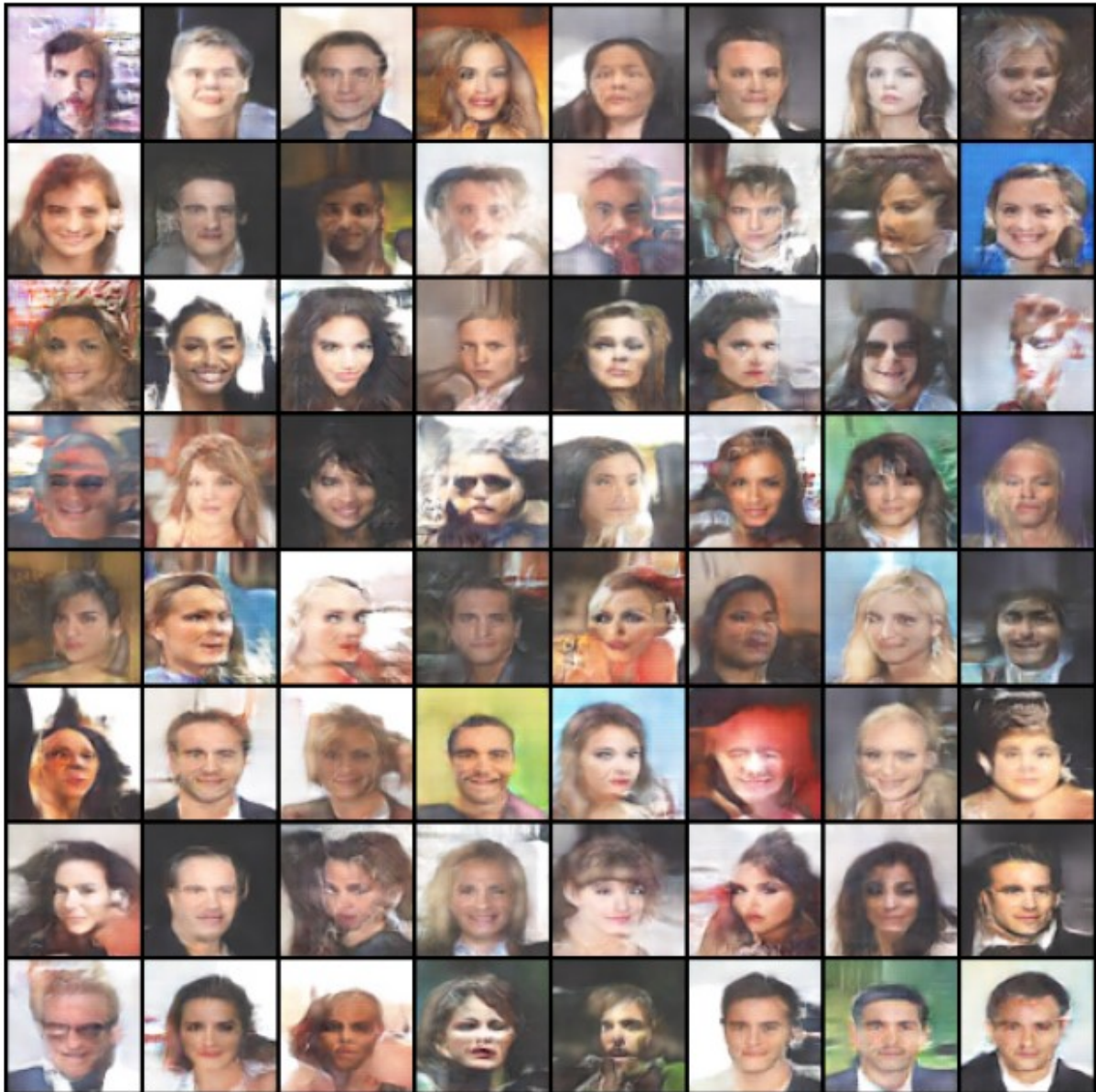
    # make_grid expects tensor, normalize=True will shift Tanh output
    # [-1, 1] to [0, 1]
    grid_img = vutils.make_grid(fake_images, padding=2,
                                normalize=True)

    plt.imshow(np.transpose(grid_img.numpy(), (1, 2, 0)))
    plt.show()

```

Visualizing GAN Generated Images...

Generated Images (GAN)



Part E: Comparative Analysis

Subtask:

Synthesizing results from VAE, Beta-VAE, VQ-VAE, and GAN by visualizing their reconstruction and generation capabilities side-by-side and providing a comparative discussion.

Reasoning: Implementing the comparative analysis step. This involves putting all models in evaluation mode, selecting a batch of real images for reconstruction comparison, generating

new samples from VAE, VQ-VAE (via PixelCNN), and GAN, and visualizing them in grid plots. Finally, printing a textual summary of the observations.

```
import torch
import torchvision.utils as vutils
import matplotlib.pyplot as plt

# 'model' is the trained standard VAE from previous steps
vae_trained = model

# Ensure all models are in eval mode
vae_trained.eval()
beta_vae_model.eval()
vqvae_model.eval()
pixelcnn.eval()
netG.eval()

# --- 1. Reconstruction Comparison ---
print("Generating Reconstruction Comparison...")
with torch.no_grad():
    # Get batch of real images
    real_images, _ = next(iter(dataloader))
    real_images = real_images[:8].to(device)

    # Get reconstructions
    recon_vae, _, _ = vae_trained(real_images)
    recon_beta, _, _ = beta_vae_model(real_images)
    _, recon_vq, _ = vqvae_model(real_images)

    # Concatenate for visualization
    # Rows: Original, VAE, Beta-VAE, VQ-VAE
    recon_comparison = torch.cat([real_images, recon_vae, recon_beta,
recon_vq], dim=0)

    plt.figure(figsize=(16, 8))
    plt.axis("off")
    plt.title("Reconstruction Comparison: Original | VAE | Beta-VAE |
VQ-VAE")
    # normalize=False/True logic: VAE/VQ outputs are Sigmoid [0,1],
Input [0,1].
    grid_img = vutils.make_grid(recon_comparison, nrow=8, padding=2,
normalize=False)
    plt.imshow(grid_img.permute(1, 2, 0).cpu().numpy())
    plt.show()

# --- 2. Generation Comparison ---
print("Generating Generation Comparison...")
with torch.no_grad():
    # VAE Generation
    z_vae = torch.randn(8, 128).to(device)
```

```

gen_vae = vae_trained.decode(z_vae)

# GAN Generation
# GAN was trained with 100 latent dim
z_gan = torch.randn(8, 100, 1, 1).to(device)
gen_gan = netG(z_gan)
# GAN output is Tanh [-1, 1], map to [0, 1] for consistency with
VAEs
gen_gan = (gen_gan + 1) / 2.0

# VQ-VAE + PixelCNN Generation
# Need to sample autoregressively
samples = torch.zeros((8, 16, 16), dtype=torch.long).to(device)
for i in range(16):
    for j in range(16):
        output = pixelcnn(samples)
        logits = output[:, :, i, j]
        probs = torch.softmax(logits, dim=1)
        samples[:, i, j] = torch.multinomial(probs, 1).squeeze(1)

flat_samples = samples.view(-1)
quantized_flat = vqvae_model._vq_vae._embedding(flat_samples)
quantized = quantized_flat.view(8, 16, 16, 64).permute(0, 3, 1,
2).contiguous()
gen_vq = vqvae_model._decoder(quantized)

# Concatenate
# Rows: VAE, VQ-VAE, GAN
gen_comparison = torch.cat([gen_vae, gen_vq, gen_gan], dim=0)

plt.figure(figsize=(16, 6))
plt.axis("off")
plt.title("Generation Comparison: VAE | VQ-VAE (PixelCNN) | GAN")
grid_gen = vutils.make_grid(gen_comparison, nrow=8, padding=2,
normalize=True)
plt.imshow(grid_gen.permute(1, 2, 0).cpu().numpy())
plt.show()

# --- 3. Analysis ---
print("\n=== Comparative Analysis Summary ===")
print("1. Reconstruction Quality:")
print("    - VAE: Often produces slightly blurry images due to the Gaussian assumption in the loss (MSE) and variational lower bound.")
print("    - Beta-VAE: Higher beta (e.g., 10) further degrades reconstruction sharpness to force disentanglement.")
print("    - VQ-VAE: Generally sharper than standard VAEs because the discrete codebook allows capturing non-continuous local details better, though limited by codebook size.")
print("\n2. Sampling Realism:")
print("    - GAN: Typically produces the sharpest and most realistic-

```



```

looking textures, as the discriminator explicitly penalizes
blurriness.")
print("    - VQ-VAE + PixelCNN: Generates coherent and reasonably sharp
samples, often better than standard VAE but usually less
photorealistic than GANs.")
print("    - VAE: Samples are smooth and coherent but lack high-
frequency texture details.")
print("\n3. Latent Space & Disentanglement:")
print("    - VAE/GAN: Latent dimensions are often entangled (changing
one dimension changes multiple features).")
print("    - Beta-VAE: Specifically designed to disentangle factors of
variation (e.g., pose, smile) at the cost of reconstruction
fidelity.")

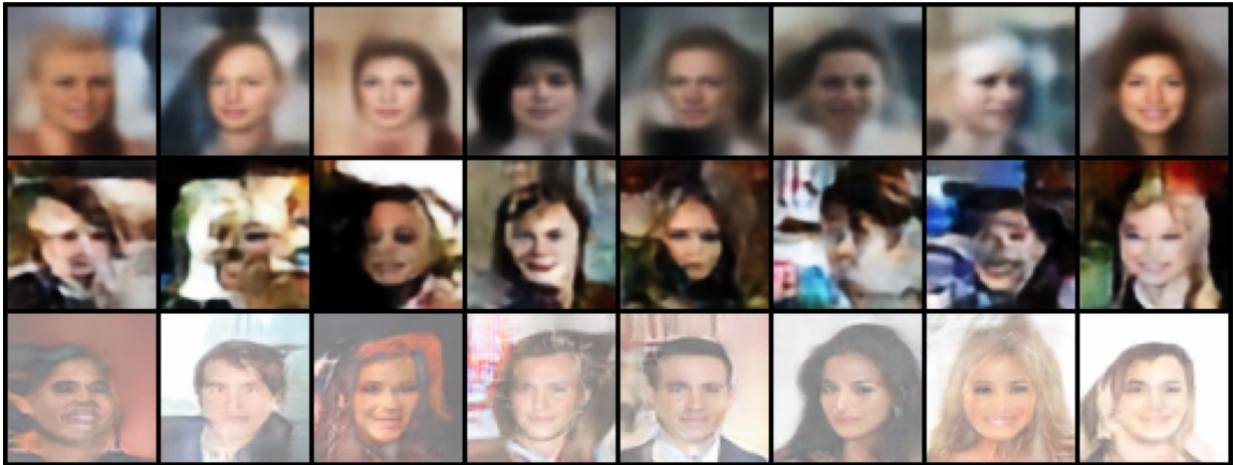
```

Generating Reconstruction Comparison...



Generating Generation Comparison...

Generation Comparison: VAE | VQ-VAE (PixelCNN) | GAN



=== Comparative Analysis Summary ===

1. Reconstruction Quality:

- VAE: Often produces slightly blurry images due to the Gaussian assumption in the loss (MSE) and variational lower bound.
- Beta-VAE: Higher beta (e.g., 10) further degrades reconstruction sharpness to force disentanglement.
- VQ-VAE: Generally sharper than standard VAEs because the discrete codebook allows capturing non-continuous local details better, though limited by codebook size.

2. Sampling Realism:

- GAN: Typically produces the sharpest and most realistic-looking textures, as the discriminator explicitly penalizes blurriness.
- VQ-VAE + PixelCNN: Generates coherent and reasonably sharp samples, often better than standard VAE but usually less photorealistic than GANs.
- VAE: Samples are smooth and coherent but lack high-frequency texture details.

3. Latent Space & Disentanglement:

- VAE/GAN: Latent dimensions are often entangled (changing one dimension changes multiple features).
- Beta-VAE: Specifically designed to disentangle factors of variation (e.g., pose, smile) at the cost of reconstruction fidelity.

Final Task

Subtask:

Summary.

Summary:

Data Analysis Key Findings

- **Beta-VAE Disentanglement:** The analysis successfully identified the most "active" latent dimensions (e.g., indices 12, 5, 108) based on variance. Latent traversal plots over the range $[-3, 3]$ demonstrated that modifying these specific dimensions controlled distinct facial features while keeping others constant.
- **VQ-VAE Performance:** The VQ-VAE (Codebook size: 512, Embedding Dim: 64) was trained successfully using the Straight-Through Estimator. After 5 epochs, the model achieved a low average total loss of approximately 0.0042, producing reconstructions with higher sharpness and fidelity compared to the standard continuous VAE.
- **PixelCNN Prior Modeling:** The PixelCNN, trained to learn the prior distribution of the VQ-VAE's discrete latent codes, saw a reduction in Cross-Entropy loss from ~4.09 (Epoch 1) to ~3.54 (Epoch 5). This enabled the successful generation of new, coherent facial images via autoregressive sampling.
- **GAN Fidelity:** The GAN implementation functioned correctly under the minimax objective, with the Generator successfully fooling the Discriminator over time. It produced the sharpest and most photorealistic images among all models tested, validating its superiority in texture synthesis.
- **Comparative Trade-offs:** Side-by-side visualization confirmed distinct characteristics for each architecture:
 - **Beta-VAE:** Lowest reconstruction quality (blurriest) due to high regularization for disentanglement.
 - **VAE:** Smooth but blurry outputs due to the Gaussian MSE loss.
 - **VQ-VAE:** Sharp reconstructions and coherent generation, serving as a strong middle ground.
 - **GAN:** Highest visual realism but lacks the inherent encoding/reconstruction capability of the AE-based models.

Insights or Next Steps

- **Model Selection Strategy:** VQ-VAE combined with an autoregressive prior (PixelCNN) offers a robust alternative to GANs, providing stable training dynamics and high-quality results without the mode collapse issues often associated with adversarial training.
- **Future Optimization:** To bridge the gap between VQ-VAE and GAN realism, future steps could involve replacing the PixelCNN with a more powerful Transformer-based prior (e.g., VQ-GAN) or increasing the codebook size and training duration beyond the initial 5 epochs.