



Creating Numbers/images with AI: A Hands-on Diffusion Model Exercise

Introduction

In this assignment, you'll learn how to create an AI model that can generate realistic images from scratch using a powerful technique called 'diffusion'. Think of it like teaching AI to draw by first learning how images get blurry and then learning to make them clear again.

What We'll Build

- A diffusion model capable of generating realistic images
- For most students: An AI that generates handwritten digits (0-9) using the MNIST dataset
- For students with more computational resources: Options to work with more complex datasets
- Visual demonstrations of how random noise gradually transforms into clear, recognizable images
- By the end, your AI should create images realistic enough for another AI to recognize them

Dataset Options

This lab offers flexibility based on your available computational resources:

- Standard Option (Free Colab): We'll primarily use the MNIST handwritten digit dataset, which works well with limited GPU memory and completes training in a reasonable time frame. Most examples and code in this notebook are optimized for MNIST.
- Advanced Option: If you have access to more powerful GPUs (either through Colab Pro/Pro+ or your own hardware), you can experiment with more complex datasets like Fashion-MNIST, CIFAR-10, or even face generation. You'll need to adapt the model architecture, hyperparameters, and evaluation metrics accordingly.

Resource Requirements

- Basic MNIST: Works with free Colab GPUs (2-4GB VRAM), ~30 minutes

training

- Fashion-MNIST: Similar requirements to MNIST CIFAR-10: Requires more memory (8-12GB VRAM) and longer training (~2 hours)
- Higher resolution images: Requires substantial GPU resources and several hours of training

Before You Start

1. Make sure you're running this in Google Colab or another environment with GPU access
2. Go to 'Runtime' → 'Change runtime type' and select 'GPU' as your hardware accelerator
3. Each code cell has comments explaining what it does
4. Don't worry if you don't understand every detail - focus on the big picture!
5. If working with larger datasets, monitor your GPU memory usage carefully

The concepts you learn with MNIST will scale to more complex datasets, so even if you're using the basic option, you'll gain valuable knowledge about generative AI that applies to more advanced applications.

Step 1: Setting Up Our Tools

First, let's install and import all the tools we need. Run this cell and wait for it to complete.

```
In [210... # Step 1: Install required packages
%pip install einops
print("Package installation complete.")

# Step 2: Import libraries
# --- Core PyTorch libraries ---
import torch # Main deep learning framework
import torch.nn.functional as F # Neural network functions like activation fu
import torch.nn as nn # Neural network building blocks (layers)
from torch.optim import Adam # Optimization algorithm for training

# --- Data handling ---
from torch.utils.data import Dataset, DataLoader # For organizing and loading
import torchvision # Library for computer vision datasets and models
import torchvision.transforms as transforms # For preprocessing images

# --- Tensor manipulation ---
import random # For random operations
```

```

from einops.layers.torch import Rearrange # For reshaping tensors in neural networks
from einops import rearrange # For elegant tensor reshaping operations
import numpy as np # For numerical operations on arrays

# --- System utilities ---
import os # For operating system interactions (used for CPU count)

# --- Visualization tools ---
import matplotlib.pyplot as plt # For plotting images and graphs
from PIL import Image # For image processing
from torchvision.utils import save_image, make_grid # For saving and displaying images

# Step 3: Set up device (GPU or CPU)
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print(f"We'll be using: {device}")

# Check if we're actually using GPU (for students to verify)
if device.type == "cuda":
    print(f"GPU name: {torch.cuda.get_device_name(0)}")
    print(f"GPU memory: {torch.cuda.get_device_properties(0).total_memory / 1e9} GB")
else:
    print("Note: Training will be much slower on CPU. Consider using Google Colab.")

```

Requirement already satisfied: einops in /usr/local/lib/python3.11/dist-packages (0.8.1)

Package installation complete.

We'll be using: cuda

GPU name: Tesla T4

GPU memory: 15.83 GB

REPRODUCIBILITY AND DEVICE SETUP

```

In [211]: # Step 4: Set random seeds for reproducibility
# Diffusion models are sensitive to initialization, so reproducible results here
SEED = 42 # Universal seed value for reproducibility
torch.manual_seed(SEED) # PyTorch random number generator
np.random.seed(SEED) # NumPy random number generator
random.seed(SEED) # Python's built-in random number generator

print(f"Random seeds set to {SEED} for reproducible results")

# Configure CUDA for GPU operations if available
if torch.cuda.is_available():
    torch.cuda.manual_seed(SEED) # GPU random number generator
    torch.cuda.manual_seed_all(SEED) # All GPUs random number generator

# Ensure deterministic GPU operations
# Note: This slightly reduces performance but ensures results are reproducible
torch.backends.cudnn.deterministic = True
torch.backends.cudnn.benchmark = False

try:
    # Check available GPU memory

```

```

gpu_memory = torch.cuda.get_device_properties(0).total_memory / 1e9 #
print(f"Available GPU Memory: {gpu_memory:.1f} GB")

# Add recommendation based on memory
if gpu_memory < 4:
    print("Warning: Low GPU memory. Consider reducing batch size if yo
except Exception as e:
    print(f"Could not check GPU memory: {e}")
else:
    print("No GPU detected. Training will be much slower on CPU.")
    print("If you're using Colab, go to Runtime > Change runtime type and sele

```

Random seeds set to 42 for reproducible results
Available GPU Memory: 15.8 GB

Step 2: Choosing Your Dataset

You have several options for this exercise, depending on your computer's capabilities:

Option 1: MNIST (Basic - Works on Free Colab)

- Content: Handwritten digits (0-9)
- Image size: 28x28 pixels, Grayscale
- Training samples: 60,000
- Memory needed: ~2GB GPU
- Training time: ~15-30 minutes on Colab
- **Choose this if:** You're using free Colab or have a basic GPU

Option 2: Fashion-MNIST (Intermediate)

- Content: Clothing items (shirts, shoes, etc.)
- Image size: 28x28 pixels, Grayscale
- Training samples: 60,000
- Memory needed: ~2GB GPU
- Training time: ~15-30 minutes on Colab
- **Choose this if:** You want more interesting images but have limited GPU

Option 3: CIFAR-10 (Advanced)

- Content: Real-world objects (cars, animals, etc.)
- Image size: 32x32 pixels, Color (RGB)
- Training samples: 50,000
- Memory needed: ~4GB GPU

- Training time: ~1-2 hours on Colab
- **Choose this if:** You have Colab Pro or a good local GPU (8GB+ memory)

Option 4: CelebA (Expert)

- Content: Celebrity face images
- Image size: 64x64 pixels, Color (RGB)
- Training samples: 200,000
- Memory needed: ~8GB GPU
- Training time: ~3-4 hours on Colab
- **Choose this if:** You have excellent GPU (12GB+ memory)

To use your chosen dataset, uncomment its section in the code below and make sure all others are commented out.

In [212...

```
#=====
# SECTION 2: DATASET SELECTION AND CONFIGURATION
#=====
# STUDENT INSTRUCTIONS:
# 1. Choose ONE dataset option based on your available GPU memory
# 2. Uncomment ONLY ONE dataset section below
# 3. Make sure all other dataset sections remain commented out

#-----
# OPTION 1: MNIST (Basic - 2GB GPU)
#-----
# Recommended for: Free Colab or basic GPU
# Memory needed: ~2GB GPU
# Training time: ~15-30 minutes

IMG_SIZE = 28
IMG_CH = 1
N_CLASSES = 10
BATCH_SIZE = 64
EPOCHS = 30

transform = transforms.Compose([
    transforms.ToTensor(),
    transforms.Normalize((0.5,), (0.5,))
])

# Your code to load the MNIST dataset
# Hint: Use torchvision.datasets.MNIST with root='./data', train=True,
#       transform=transform, and download=True
# Then print a success message

# Enter your code her# Load MNIST dataset
dataset = torchvision.datasets.MNIST(
```

```

    root='./data',
    train=True,
    transform=transform,
    download=True
)

print("💎 MNIST dataset loaded successfully!")
print(f"Dataset size: {len(dataset)} images")
print(f"Image size: {IMG_SIZE}x{IMG_SIZE}")
print(f"Number of classes: {N_CLASSES}")
print(f"Batch size: {BATCH_SIZE}")
print(f"Training epochs: {EPOCHS}")

#-----
# OPTION 2: Fashion-MNIST (Intermediate - 2GB GPU)
#-----
# Uncomment this section to use Fashion-MNIST instead
"""
IMG_SIZE = 28
IMG_CH = 1
N_CLASSES = 10
BATCH_SIZE = 64
EPOCHS = 30

transform = transforms.Compose([
    transforms.ToTensor(),
    transforms.Normalize((0.5,), (0.5,))
])

# Your code to load the Fashion-MNIST dataset
# Hint: Very similar to MNIST but use torchvision.datasets.FashionMNIST

# Enter your code here:

"""

#-----
# OPTION 3: CIFAR-10 (Advanced - 4GB+ GPU)
#-----
# Uncomment this section to use CIFAR-10 instead
"""
IMG_SIZE = 32
IMG_CH = 3
N_CLASSES = 10
BATCH_SIZE = 32 # Reduced batch size for memory
EPOCHS = 50     # More epochs for complex data

# Your code to create the transform and load CIFAR-10
# Hint: Use transforms.Normalize with RGB means and stds ((0.5, 0.5, 0.5), (0.
# Then load torchvision.datasets.CIFAR10

```

```
# Enter your code here:
```

```
"""
```

❖ MNIST dataset loaded successfully!

Dataset size: 60000 images

Image size: 28x28

Number of classes: 10

Batch size: 64

Training epochs: 30

```
Out[212...] '\nIMG_SIZE = 32\nIMG_CH = 3\nN_CLASSES = 10\nBATCH_SIZE = 32 # Reduced batch size for memory\nEPOCHS = 50 # More epochs for complex data\n\n# Your code to create the transform and load CIFAR-10\n# Hint: Use transforms.Normalize with RGB means and stds ((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))\n# Then load torchvision.datasets.CIFAR10\n# Enter your code here:\n\n'
```

```
In [213...] #Validating Dataset Selection
#Let's add code to validate that a dataset was selected
# and check if your GPU has enough memory:

# Validate dataset selection
if 'dataset' not in locals():
    raise ValueError("""
❖ ERROR: No dataset selected! Please uncomment exactly one dataset option
Available options:
1. MNIST (Basic) - 2GB GPU
2. Fashion-MNIST (Intermediate) - 2GB GPU
3. CIFAR-10 (Advanced) - 4GB+ GPU
4. CelebA (Expert) - 8GB+ GPU
""")

# Validate GPU memory requirements
if torch.cuda.is_available():
    gpu_memory = torch.cuda.get_device_properties(0).total_memory / 1e9 # Convert to GB
    print(f"Available GPU Memory: {gpu_memory:.1f} GB")

    if gpu_memory < 2:
        print("Warning: Low GPU memory. Consider reducing batch size if you encounter errors.")
    else:
        print("No GPU detected. Training will be much slower on CPU.")
```

Available GPU Memory: 15.8 GB

```
In [214...] #Dataset Properties and Data Loaders
#Now let's examine our dataset
#and set up the data loaders:

# Your code to check sample batch properties
# Hint: Get a sample batch using next(iter(DataLoader(dataset, batch_size=1)))
# Then print information about the dataset shape, type, and value ranges

# Enter your code here:
sample_dataloader = DataLoader(dataset, batch_size=1, shuffle=True)
```

```

sample_batch = next(iter(sample_dataloader))
sample_images, sample_labels = sample_batch

print("Dataset Properties:")
print(f"Sample image shape: {sample_images.shape}")
print(f"Sample label shape: {sample_labels.shape}")
print(f"Image data type: {sample_images.dtype}")
print(f"Label data type: {sample_labels.dtype}")
print(f"Image value range: [{sample_images.min().item():.3f}, {sample_images.max().item():.3f}]")
print(f"Sample label: {sample_labels.item()}")
print(f"Total dataset size: {len(dataset)} images")

#=====
# SECTION 3: DATASET SPLITTING AND DATALOADER CONFIGURATION
#=====
# Create train-validation split

# Your code to create a train-validation split (80% train, 20% validation)
# Hint: Use random_split() with appropriate train_size and val_size
# Be sure to use a fixed generator for reproducibility

# Enter your code here:
total_size = len(dataset)
train_size = int(0.8 * total_size) # 80% for training
val_size = total_size - train_size # 20% for validation

print(f"Dataset splitting:")
print(f"Total images: {total_size}")
print(f"Training images: {train_size}")
print(f"Validation images: {val_size}")

# Create train-validation split with fixed generator for reproducibility
train_dataset, val_dataset = torch.utils.data.random_split(
    dataset,
    [train_size, val_size],
    generator=torch.Generator().manual_seed(SEED)
)

print("💎 Train-validation split created successfully!")

# Your code to create dataloaders for training and validation
# Hint: Use DataLoader with batch_size=BATCH_SIZE, appropriate shuffle setting
# and num_workers based on available CPU cores

# Enter your code here:
# Determine number of workers based on available CPU cores
num_workers = min(4, os.cpu_count()) # Use up to 4 workers or available cores

# Create training dataloader
train_dataloader = DataLoader(
    train_dataset,
    batch_size=BATCH_SIZE,

```



```

        shuffle=True, # Shuffle training data
        num_workers=num_workers,
        pin_memory=True if torch.cuda.is_available() else False # Faster data tra
    )

# Create validation dataloader
val_dataloader = DataLoader(
    val_dataset,
    batch_size=BATCH_SIZE,
    shuffle=False, # No need to shuffle validation data
    num_workers=num_workers,
    pin_memory=True if torch.cuda.is_available() else False
)

print(f"💎 DataLoaders created successfully!")
print(f"Training batches: {len(train_dataloader)}")
print(f"Validation batches: {len(val_dataloader)}")
print(f"Number of workers: {num_workers}")
print(f"GPU acceleration: {'Enabled' if torch.cuda.is_available() else 'Disabl

# Verify the dataloaders work correctly
print("\nVerifying dataloaders...")
try:
    # Test training dataloader
    train_batch = next(iter(train_dataloader))
    train_images, train_labels = train_batch
    print(f"Training batch shape: {train_images.shape}")
    print(f"Training labels shape: {train_labels.shape}")

    # Test validation dataloader
    val_batch = next(iter(val_dataloader))
    val_images, val_labels = val_batch
    print(f"Validation batch shape: {val_images.shape}")
    print(f"Validation labels shape: {val_labels.shape}")

    print("💎 DataLoaders working correctly!")

except Exception as e:
    print(f"💎 Error with dataloaders: {e}")

```

```

Dataset Properties:
Sample image shape: torch.Size([1, 1, 28, 28])
Sample label shape: torch.Size([1])
Image data type: torch.float32
Label data type: torch.int64
Image value range: [-1.000, 0.992]
Sample label: 1
Total dataset size: 60000 images
Dataset splitting:
Total images: 60000
Training images: 48000
Validation images: 12000
✦ Train-validation split created successfully!
✦ DataLoaders created successfully!
Training batches: 750
Validation batches: 188
Number of workers: 2
GPU acceleration: Enabled

```

```

Verifying dataloaders...
Training batch shape: torch.Size([64, 1, 28, 28])
Training labels shape: torch.Size([64])
Validation batch shape: torch.Size([64, 1, 28, 28])
Validation labels shape: torch.Size([64])
✦ DataLoaders working correctly!

```

Step 3: Building Our Model Components

Now we'll create the building blocks of our AI model. Think of these like LEGO pieces that we'll put together to make our number generator:

- GELUConvBlock: The basic building block that processes images
- DownBlock: Makes images smaller while finding important features
- UpBlock: Makes images bigger again while keeping the important features
- Other blocks: Help the model understand time and what number to generate

```

In [215... # Basic building block that processes images
class GELUConvBlock(nn.Module):
    def __init__(self, in_ch, out_ch, group_size):
        """
        Creates a block with convolution, normalization, and activation

        Args:
            in_ch (int): Number of input channels
            out_ch (int): Number of output channels
            group_size (int): Number of groups for GroupNorm
        """
        super().__init__()

```

```

# Check that group_size is compatible with out_ch
if out_ch % group_size != 0:
    print(f"Warning: out_ch ({out_ch}) is not divisible by group_size")
    # Adjust group_size to be compatible
    group_size = min(group_size, out_ch)
    while out_ch % group_size != 0:
        group_size -= 1
    print(f"Adjusted group_size to {group_size}")

# Your code to create layers for the block
# Hint: Use nn.Conv2d, nn.GroupNorm, and nn.GELU activation
# Then combine them using nn.Sequential

# Enter your code here:
self.conv_block = nn.Sequential(
    nn.Conv2d(in_ch, out_ch, kernel_size=3, padding=1),
    nn.GroupNorm(group_size, out_ch),
    nn.GELU()
)

def forward(self, x):
    # Your code for the forward pass
    # Hint: Simply pass the input through the model

    # Enter your code here:
    return self.conv_block(x)

pass

```

```

In [216... # Rearranges pixels to downsample the image (2x reduction in spatial dimension)
class RearrangePoolBlock(nn.Module):
    def __init__(self, in_chs, group_size):
        """
        Downsamples the spatial dimensions by 2x while preserving information

        Args:
            in_chs (int): Number of input channels
            group_size (int): Number of groups for GroupNorm
        """
        super().__init__()

        # Your code to create the rearrange operation and convolution
        # Hint: Use Rearrange from einops.layers.torch to reshape pixels
        # Then add a GELUConvBlock to process the rearranged tensor

        # Enter your code here:
        # Rearrange operation to downsample by 2x
        self.rearrange = Rearrange('b c (h p1) (w p2) -> b (c p1 p2) h w', p1=

        # Convolution block to process the rearranged tensor
        self.conv_block = GELUConvBlock(in_chs * 4, in_chs, group_size)

```

```

def forward(self, x):
    # Your code for the forward pass
    # Hint: Apply rearrange to downsample, then apply convolution

    # Enter your code here:
    x = self.rearrange(x)
    x = self.conv_block(x)
    return x
pass

```

```

In [217... #Let's implement the upsampling block for our U-Net architecture:
class DownBlock(nn.Module):
    """
    Downsampling block for encoding path in U-Net architecture.

    This block:
    1. Processes input features with two convolutional blocks
    2. Downsamples spatial dimensions by 2x using pixel rearrangement

    Args:
        in_chs (int): Number of input channels
        out_chs (int): Number of output channels
        group_size (int): Number of groups for GroupNorm
    """
    def __init__(self, in_chs, out_chs, group_size):
        super().__init__() # Simplified super() call, equivalent to original
        # Removed time and class projections from __init__
        # self.time_proj = nn.Conv2d(time_dim, ch[-1], 1) # Project to match
        # self.class_proj = nn.Conv2d(class_dim, ch[-1], 1)
        # Sequential processing of features
        layers = [
            GELUConvBlock(in_chs, out_chs, group_size), # First conv block ch
            GELUConvBlock(out_chs, out_chs, group_size), # Second conv block
            RearrangePoolBlock(out_chs, group_size) # Downsampling (spat
        ]
        self.model = nn.Sequential(*layers)

        # Log the configuration for debugging
        print(f"Created DownBlock: in_chs={in_chs}, out_chs={out_chs}, spatial

    def forward(self, x):
        # Removed time and class embedding processing from forward
        # t_emb = self.time_embed(t) # Shape: [batch_size, time_dim]
        # t_emb = t_emb.view(t_emb.shape[0], t_emb.shape[1], 1, 1) # Shape: [
        # t_emb = self.time_proj(t_emb) # Project to match feature map channe

        # c_emb = self.class_embed(c) # Shape: [batch_size, class_dim]
        # c_emb = c_emb.view(c_emb.shape[0], c_emb.shape[1], 1, 1) # Shape: [
        # c_emb = self.class_proj(c_emb) # Project to match feature map channr
        # c_emb = c_emb * c_mask.view(-1, 1, 1, 1) # Apply class mask
        return self.model(x)

```

```

In [218... #Now let's implement the upsampling block for our U-Net architecture:

```

```

class UpBlock(nn.Module):
    """
    Upsampling block for decoding path in U-Net architecture.

    This block:
    1. Takes features from the decoding path and corresponding skip connection
    2. Concatenates them along the channel dimension
    3. Upsamples spatial dimensions by 2x using transposed convolution
    4. Processes features through multiple convolutional blocks

    Args:
        in_chs (int): Number of input channels from the previous layer
        out_chs (int): Number of output channels
        group_size (int): Number of groups for GroupNorm
    """
    def __init__(self, in_chs, out_chs, group_size):
        super().__init__()

        # Your code to create the upsampling operation
        # Hint: Use nn.ConvTranspose2d with kernel_size=2 and stride=2
        # Note that the input channels will be 2 * in_chs due to concatenation

        # Enter your code here:
        # Upsampling operation using transposed convolution
        self.upsample = nn.ConvTranspose2d(
            in_channels=2 * in_chs, # 2x due to concatenation with skip connection
            out_channels=2 * in_chs, # Keep same number of channels after upsampling
            kernel_size=2,
            stride=2
        )

        # Convolutional blocks for feature processing
        layers = [
            GELUConvBlock(2 * in_chs, out_chs, group_size), # First conv block
            GELUConvBlock(out_chs, out_chs, group_size),    # Second conv block
        ]
        self.conv_blocks = nn.Sequential(*layers)

        # Log the configuration for debugging
        print(f"Created UpBlock: in_chs={in_chs}, out_chs={out_chs}, spatial_i

    def forward(self, x, skip):
        """
        Forward pass through the UpBlock.

        Args:
            x (torch.Tensor): Input tensor from previous layer [B, in_chs, H, W]
            skip (torch.Tensor): Skip connection tensor from encoder [B, in_chs, H, W]

        Returns:
            torch.Tensor: Output tensor with shape [B, out_chs, 2H, 2W]
        """

```

```

# Your code for the forward pass
# Hint: Concatenate x and skip, then upsample and process

# Enter your code here:
# Concatenate input features with skip connection along channel dimension
x = torch.cat([x, skip], dim=1) # [B, 2*in_chs, H, W]

# Apply upsampling using transposed convolution
x = self.upsample(x) # [B, 2*in_chs, 2H, 2W]

# Process features through convolutional blocks
x = self.conv_blocks(x) # [B, out_chs, 2H, 2W]

return x

```

```

In [219]: # Here we implement the time embedding block for our U-Net architecture:
# Helps the model understand time steps in diffusion process
class SinusoidalPositionEmbedBlock(nn.Module):
    """
    Creates sinusoidal embeddings for time steps in diffusion process.

    This embedding scheme is adapted from the Transformer architecture and
    provides a unique representation for each time step that preserves
    relative distance information.

    Args:
        dim (int): Embedding dimension
    """
    def __init__(self, dim):
        super().__init__()
        self.dim = dim

    def forward(self, time):
        """
        Computes sinusoidal embeddings for given time steps.

        Args:
            time (torch.Tensor): Time steps tensor of shape [batch_size]

        Returns:
            torch.Tensor: Time embeddings of shape [batch_size, dim]
        """
        device = time.device
        half_dim = self.dim // 2
        embeddings = torch.log(torch.tensor(10000.0, device=device)) / (half_dim - 1)
        embeddings = torch.exp(torch.arange(half_dim, device=device) * -embeddings)
        embeddings = time[:, None] * embeddings[None, :]
        embeddings = torch.cat((embeddings.sin(), embeddings.cos()), dim=-1)
        return embeddings

```

```

In [220]: # Helps the model understand which number/image to draw (class conditioning)
class EmbedBlock(nn.Module):
    """

```

Creates embeddings for class conditioning in diffusion models.

This module transforms a one-hot or index representation of a class into a rich embedding that can be added to feature maps.

Args:

input_dim (int): Input dimension (typically number of classes)
emb_dim (int): Output embedding dimension

"""

```
def __init__(self, input_dim, emb_dim):  
    super(EmbedBlock, self).__init__()  
    self.input_dim = input_dim
```

Your code to create the embedding layers

Hint: Use nn.Linear layers with a GELU activation, followed by

nn.Unflatten to reshape for broadcasting with feature maps

Enter your code here:

```
self.model = nn.Sequential(  
    nn.Linear(input_dim, emb_dim),  
    nn.GELU(),  
    nn.Linear(emb_dim, emb_dim),  
    nn.GELU(),  
    nn.Linear(emb_dim, emb_dim),  
    nn.Unflatten(1, (emb_dim, 1, 1))  
)
```

```
def forward(self, x):
```

"""

Computes class embeddings for the given class indices.

Args:

x (torch.Tensor): Class indices or one-hot encodings [batch_size,

Returns:

torch.Tensor: Class embeddings of shape [batch_size, emb_dim, 1, 1
(ready to be added to feature maps)

"""

```
x = x.view(-1, self.input_dim)  
return self.model(x)
```

In [250... *# Main U-Net model that puts everything together*

```
class UNet(nn.Module):
```

"""

U-Net architecture for diffusion models with time and class conditioning.

This architecture follows the standard U-Net design with:

1. Downsampling path that reduces spatial dimensions
2. Middle processing blocks
3. Upsampling path that reconstructs spatial dimensions
4. Skip connections between symmetric layers

The model is conditioned on:

- Time step (where we are in the diffusion process)
- Class labels (what we want to generate)

Args:

`T (int)`: Number of diffusion time steps
`img_ch (int)`: Number of image channels
`img_size (int)`: Size of input images
`down_chs (list)`: Channel dimensions for each level of U-Net
`t_embed_dim (int)`: Dimension for time embeddings
`c_embed_dim (int)`: Dimension for class embeddings
`group_size (int)`: Number of groups for GroupNorm layers

"""

```
def __init__(self, T, img_ch, img_size, down_chs, t_embed_dim, c_embed_dim, group_size):
    super().__init__()

    # Store key parameters
    self.T = T
    self.img_ch = img_ch
    self.img_size = img_size
    self.down_chs = down_chs
    self.t_embed_dim = t_embed_dim
    self.c_embed_dim = c_embed_dim
    self.group_size = group_size # Store group size

    # Your code to create the time embedding
    # Hint: Use SinusoidalPositionEmbedBlock, nn.Linear, and nn.GELU in self
    self.time_embed = nn.Sequential(
        SinusoidalPositionEmbedBlock(t_embed_dim),
        nn.Linear(t_embed_dim, t_embed_dim),
        nn.GELU(),
        nn.Linear(t_embed_dim, t_embed_dim)
    )

    # Your code to create the class embedding
    # Hint: Use the EmbedBlock class you defined earlier
    self.class_embed = EmbedBlock(c_embed_dim, c_embed_dim)

    # Your code to create the initial convolution
    # Hint: Use GELUConvBlock to process the input image
    self.initial_conv = GELUConvBlock(img_ch, down_chs[0], group_size=self.group_size)

    # Your code to create the downsampling path
    # Hint: Use nn.ModuleList with DownBlock for each level
    self.down_blocks = nn.ModuleList()
    for i in range(len(down_chs) - 1):
        self.down_blocks.append(DownBlock(down_chs[i], down_chs[i + 1], group_size=self.group_size))

    # Your code to create the middle blocks
    # Hint: Use GELUConvBlock twice to process features at lowest resolution
    self.middle_blocks = nn.Sequential(
        GELUConvBlock(down_chs[-1], down_chs[-1], group_size=self.group_size),
        # ... (other middle blocks) ...
    )
```



```

        GELUConvBlock(down_chs[-1], down_chs[-1], group_size=self.group_size)
    )

    # Your code to create the upsampling path with skip connections
    # Hint: Use nn.ModuleList with UpBlock for each level (in reverse order)
    self.up_blocks = nn.ModuleList()
    for i in range(len(down_chs) - 1, 0, -1):
        self.up_blocks.append(UpBlock(down_chs[i], down_chs[i - 1], group_size=self.group_size))

    # Your code to create the final convolution
    # Hint: Use nn.Conv2d to project back to the original image channels
    self.final_conv = nn.Conv2d(down_chs[0], img_ch, kernel_size=3, padding=1)

    # Add projection layers for time and class embeddings to match middle block channels
    self.time_proj = nn.Linear(t_embed_dim, down_chs[-1])
    self.class_proj = nn.Linear(c_embed_dim, down_chs[-1])

    print(f"Created UNet with {len(down_chs)} scale levels")
    print(f"Channel dimensions: {down_chs}")

def forward(self, x, t, c, c_mask):
    """
    Forward pass through the UNet.

    Args:
        x (torch.Tensor): Input noisy image [B, img_ch, H, W]
        t (torch.Tensor): Diffusion time steps [B]
        c (torch.Tensor): Class labels [B, c_embed_dim]
        c_mask (torch.Tensor): Mask for conditional generation [B, 1]

    Returns:
        torch.Tensor: Predicted noise in the input image [B, img_ch, H, W]
    """
    # Your code for the time embedding
    # Hint: Process the time steps through the time embedding module
    t_emb = self.time_embed(t)
    # Project time embedding to match middle block channels
    t_emb = self.time_proj(t_emb)

    # Your code for the class embedding
    # Hint: Process the class labels through the class embedding module
    c_emb = self.class_embed(c)
    # Project class embedding to match middle block channels
    c_emb = self.class_proj(c_emb)

    # Your code for the initial feature extraction
    # Hint: Apply initial convolution to the input
    x = self.initial_conv(x)

    # Your code for the downsampling path and skip connections

```

```

# Hint: Process the features through each downsampling block
# and store the outputs for skip connections
skip_connections = []
for down_block in self.down_blocks:
    x = down_block(x)
    skip_connections.append(x)

# Your code for the middle processing and conditioning
# Hint: Process features through middle blocks, then add time and class embeddings
x = self.middle_blocks(x)
# Add projected time and class embeddings (reshaped for broadcasting)
x = x + t_emb.view(t_emb.shape[0], t_emb.shape[1], 1, 1)
x = x + c_emb.view(c_emb.shape[0], c_emb.shape[1], 1, 1)

# Your code for the upsampling path with skip connections
# Hint: Process features through each upsampling block,
# combining with corresponding skip connections
for i, up_block in enumerate(self.up_blocks):
    skip_connection = skip_connections[-(i + 1)]
    x = up_block(x, skip_connection)

# Your code for the final projection
# Hint: Apply the final convolution to get output in image space
x = self.final_conv(x)

return x

```

Step 4: Setting Up The Diffusion Process

Now we'll create the process of adding and removing noise from images. Think of it like:

1. Adding fog: Slowly making the image more and more blurry until you can't see it
2. Removing fog: Teaching the AI to gradually make the image clearer
3. Controlling the process: Making sure we can generate specific numbers we want

```

In [251... # Set up the noise schedule
n_steps = 100 # How many steps to go from clear image to noise
beta_start = 0.0001 # Starting noise level (small)
beta_end = 0.02 # Ending noise level (larger)

# Create schedule of gradually increasing noise levels
beta = torch.linspace(beta_start, beta_end, n_steps).to(device)

# Calculate important values used in diffusion equations
alpha = 1 - beta # Portion of original image to keep at each step
alpha_bar = torch.cumprod(alpha, dim=0) # Cumulative product of alphas

```

```

sqrt_alpha_bar = torch.sqrt(alpha_bar) # For scaling the original image
sqrt_one_minus_alpha_bar = torch.sqrt(1 - alpha_bar) # For scaling the noise

```

```

In [252... # Function to add noise to images (forward diffusion process)
def add_noise(x_0, t):
    """
    Add noise to images according to the forward diffusion process.

    The formula is:  $x_t = \sqrt{\alpha_{\text{bar}_t}} * x_0 + \sqrt{(1-\alpha_{\text{bar}_t})} * \epsilon$ 
    where  $\epsilon$  is random noise and  $\alpha_{\text{bar}_t}$  is the cumulative product of  $(1-\beta)$ .

    Args:
        x_0 (torch.Tensor): Original clean image [B, C, H, W]
        t (torch.Tensor): Timestep indices indicating noise level [B]

    Returns:
        tuple: (noisy_image, noise_added)
            - noisy_image is the image with noise added
            - noise_added is the actual noise that was added (for training)
    """
    # Create random Gaussian noise with same shape as image
    noise = torch.randn_like(x_0)

    # Get noise schedule values for the specified timesteps
    # Reshape to allow broadcasting with image dimensions
    sqrt_alpha_bar_t = sqrt_alpha_bar[t].reshape(-1, 1, 1, 1)
    sqrt_one_minus_alpha_bar_t = sqrt_one_minus_alpha_bar[t].reshape(-1, 1, 1, 1)

    # Apply the forward diffusion equation:
    # Mixture of original image (scaled down) and noise (scaled up) # Your
    # Hint: Mix the original image and noise according to the noise schedule

    # Enter your code here:
    x_t = sqrt_alpha_bar_t * x_0 + sqrt_one_minus_alpha_bar_t * noise

    return x_t, noise

```

```

In [253... # Function to remove noise from images (reverse diffusion process)
@torch.no_grad() # Don't track gradients during sampling (inference only)
def remove_noise(x_t, t, model, c, c_mask):
    """
    Remove noise from images using the learned reverse diffusion process.

    This implements a single step of the reverse diffusion sampling process.
    The model predicts the noise in the image, which we then use to partially
    denoise the image.

    Args:
        x_t (torch.Tensor): Noisy image at timestep t [B, C, H, W]
        t (torch.Tensor): Current timestep indices [B]

```

```

        model (nn.Module): U-Net model that predicts noise
        c (torch.Tensor): Class conditioning (what digit to generate) [B, C]
        c_mask (torch.Tensor): Mask for conditional generation [B, 1]

Returns:
    torch.Tensor: Less noisy image for the next timestep [B, C, H, W]
"""
# Predict the noise in the image using our model
predicted_noise = model(x_t, t, c, c_mask)

# Get noise schedule values for the current timestep
alpha_t = alpha[t].reshape(-1, 1, 1, 1)
alpha_bar_t = alpha_bar[t].reshape(-1, 1, 1, 1)
beta_t = beta[t].reshape(-1, 1, 1, 1)

# Special case: if we're at the first timestep (t=0), we're done
if t[0] == 0:
    return x_t
else:
    # Calculate the mean of the denoised distribution
    # This is derived from Bayes' rule and the diffusion process equations
    mean = (1 / torch.sqrt(alpha_t)) * (
        x_t - (beta_t / sqrt_one_minus_alpha_bar_t) * predicted_noise
    )

    # Add a small amount of random noise (variance depends on timestep)
    # This helps prevent the generation from becoming too deterministic
    noise = torch.randn_like(x_t)

    # Return the partially denoised image with a bit of new random noise
    return mean + torch.sqrt(beta_t) * noise

```

```

In [254... # Visualization function to show how noise progressively affects images
def show_noise_progression(image, num_steps=5):
    """
    Visualize how an image gets progressively noisier in the diffusion process

    Args:
        image (torch.Tensor): Original clean image [C, H, W]
        num_steps (int): Number of noise levels to show
    """
    plt.figure(figsize=(15, 3))

    # Show original image
    plt.subplot(1, num_steps, 1)
    if IMG_CH == 1: # Grayscale image
        plt.imshow(image[0].cpu(), cmap='gray')
    else: # Color image
        img = image.permute(1, 2, 0).cpu() # Change from [C,H,W] to [H,W,C]
        if img.min() < 0: # If normalized between -1 and 1
            img = (img + 1) / 2 # Rescale to [0,1] for display
        plt.imshow(img)
    plt.title('Original')

```

```

plt.axis('off')

# Show progressively noisier versions
for i in range(1, num_steps):
    # Calculate timestep index based on percentage through the process
    t_idx = int((i/num_steps) * n_steps)
    t = torch.tensor([t_idx]).to(device)

    # Add noise corresponding to timestep t
    noisy_image, _ = add_noise(image.unsqueeze(0), t)

    # Display the noisy image
    plt.subplot(1, num_steps, i+1)
    if IMG_CH == 1:
        plt.imshow(noisy_image[0][0].cpu(), cmap='gray')
    else:
        img = noisy_image[0].permute(1, 2, 0).cpu()
        if img.min() < 0:
            img = (img + 1) / 2
        plt.imshow(img)
    plt.title(f'{int((i/num_steps) * 100)}% Noise')
    plt.axis('off')
plt.show()

# Show an example of noise progression on a real image
sample_batch = next(iter(train_dataloader)) # Get first batch
sample_image = sample_batch[0][0].to(device) # Get first image
show_noise_progression(sample_image)

# Student Activity: Try different noise schedules
# Uncomment and modify these lines to experiment:
"""
# Try a non-linear noise schedule
beta_alt = torch.linspace(beta_start, beta_end, n_steps)**2
alpha_alt = 1 - beta_alt
alpha_bar_alt = torch.cumprod(alpha_alt, dim=0)
# How would this affect the diffusion process?
"""

# Student Activity: Try different noise schedules
# Uncomment and modify these lines to experiment:

# Try a non-linear noise schedule
beta_alt = torch.linspace(beta_start, beta_end, n_steps)**2
alpha_alt = 1 - beta_alt
alpha_bar_alt = torch.cumprod(alpha_alt, dim=0)

# Let's visualize how this alternative schedule compares to the original
plt.figure(figsize=(12, 4))

# Plot original linear schedule
plt.subplot(1, 3, 1)
plt.plot(beta.cpu(), label='Original (Linear)')
plt.plot(beta_alt, label='Alternative (Quadratic)')

```

```

plt.title('Beta Schedule Comparison')
plt.xlabel('Timestep')
plt.ylabel('Beta Value')
plt.legend()

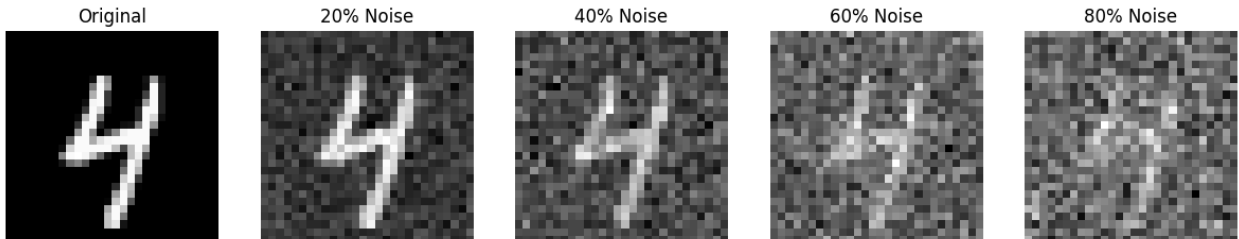
# Plot alpha_bar comparison
plt.subplot(1, 3, 2)
plt.plot(alpha_bar.cpu(), label='Original')
plt.plot(alpha_bar_alt, label='Alternative')
plt.title('Alpha Bar Schedule Comparison')
plt.xlabel('Timestep')
plt.ylabel('Alpha Bar Value')
plt.legend()

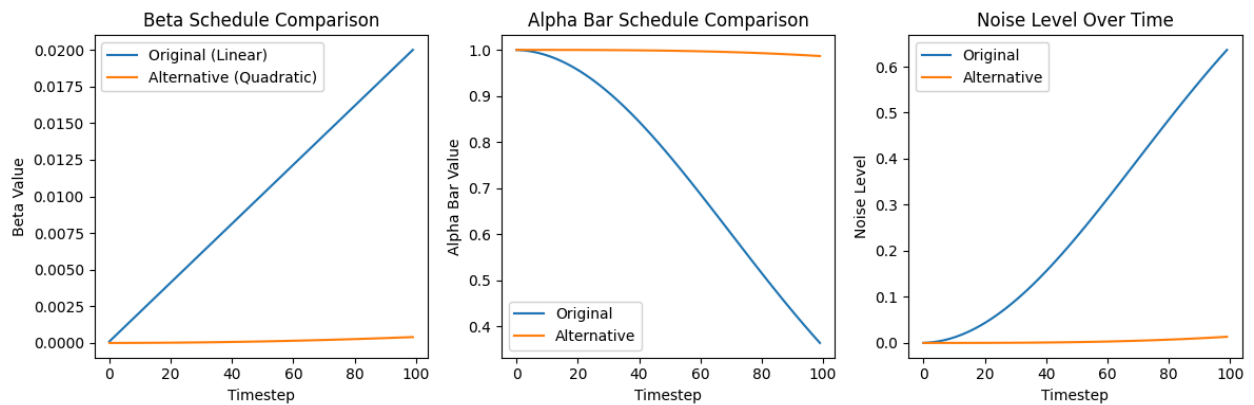
# Plot noise level over time
plt.subplot(1, 3, 3)
noise_level_orig = 1 - alpha_bar
noise_level_alt = 1 - alpha_bar_alt
plt.plot(noise_level_orig.cpu(), label='Original')
plt.plot(noise_level_alt, label='Alternative')
plt.title('Noise Level Over Time')
plt.xlabel('Timestep')
plt.ylabel('Noise Level')
plt.legend()

plt.tight_layout()
plt.show()

# How would this affect the diffusion process?
print("Analysis of the quadratic noise schedule:")
print("- The quadratic schedule adds noise more slowly at the beginning")
print("- This gives the model more time to learn fine details early in the process")
print("- However, it may make the final denoising steps more challenging")
print("- The original linear schedule provides more uniform noise addition")

```





Analysis of the quadratic noise schedule:

- The quadratic schedule adds noise more slowly at the beginning
- This gives the model more time to learn fine details early in the process
- However, it may make the final denoising steps more challenging
- The original linear schedule provides more uniform noise addition

Step 5: Training Our Model

Now we'll teach our AI to generate images. This process:

1. Takes a clear image
2. Adds random noise to it
3. Asks our AI to predict what noise was added
4. Helps our AI learn from its mistakes

This will take a while, but we'll see progress as it learns!

```
In [255... # Create our model and move it to GPU if available

model = UNet(
    n_steps,                # T: Number of diffusion time steps
    IMG_CH,                 # img_ch: Number of channels in our images (1 for gray
    IMG_SIZE,               # img_size: Size of input images (28 for MNIST, 32 for C
    (32, 64, 128),          # down_chs: Channel dimensions for each downsampling lev
    8,                      # t_embed_dim: Dimension for time step embeddings
    N_CLASSES,              # c_embed_dim: Number of classes for conditioning
    16,                     # group_size: Number of groups for GroupNorm layers
).to(device)

# Print model summary
print(f"\n{'='*50}")
print(f"MODEL ARCHITECTURE SUMMARY")
print(f"{'='*50}")
print(f"Input resolution: {IMG_SIZE}x{IMG_SIZE}")
print(f"Input channels: {IMG_CH}")
print(f"Time steps: {n_steps}")
print(f"Condition classes: {N_CLASSES}")
print(f"GPU acceleration: {'Yes' if device.type == 'cuda' else 'No'})
```

```

# Validate model parameters and estimate memory requirements
# Hint: Create functions to count parameters and estimate memory usage

# Enter your code here:

# Your code to verify data ranges and integrity
# Hint: Create functions to check data ranges in training and validation data

# Enter your code here:
def count_parameters(model):
    """Count total and trainable parameters in the model"""
    total_params = sum(p.numel() for p in model.parameters())
    trainable_params = sum(p.numel() for p in model.parameters() if p.requires_grad)
    return total_params, trainable_params

def estimate_memory_usage(model, batch_size=32):
    """Estimate memory usage for model parameters and gradients"""
    total_params, trainable_params = count_parameters(model)

    # Memory for parameters (float32 = 4 bytes)
    param_memory_mb = total_params * 4 / (1024 ** 2)

    # Memory for gradients (float32 = 4 bytes)
    grad_memory_mb = trainable_params * 4 / (1024 ** 2)

    # Rough estimate for activations (very approximate)
    # Assuming activations scale with batch size and image size
    activation_memory_mb = batch_size * IMG_SIZE * IMG_SIZE * IMG_CH * 4 / (1024 ** 2)

    total_memory_mb = param_memory_mb + grad_memory_mb + activation_memory_mb

    return {
        'parameters_mb': param_memory_mb,
        'gradients_mb': grad_memory_mb,
        'activations_mb': activation_memory_mb,
        'total_mb': total_memory_mb
    }

def check_data_ranges(dataloader, name="dataset"):
    """Check data ranges and integrity in a dataloader"""
    min_val = float('-inf')
    max_val = float('inf')
    total_samples = 0

    for batch_idx, (images, labels) in enumerate(dataloader):
        batch_min = images.min().item()
        batch_max = images.max().item()

        min_val = min(min_val, batch_min)
        max_val = max(max_val, batch_max)
        total_samples += images.size(0)

```



```

        # Check for NaN or infinite values
        if torch.isnan(images).any() or torch.isinf(images).any():
            print(f"WARNING: Found NaN or infinite values in {name} batch {bat

    # Only check first few batches for efficiency
    if batch_idx >= 5:
        break

    print(f"\n{name} Data Analysis:")
    print(f"  Samples checked: {total_samples}")
    print(f"  Value range: [{min_val:.4f}, {max_val:.4f}]")
    print(f"  Expected range: [0.0, 1.0] (normalized)")

    # Validate normalization
    if min_val < -0.1 or max_val > 1.1:
        print(f"  WARNING: Values outside expected range!")
    else:
        print(f"  ✓ Data appears properly normalized")

# Execute the validation functions
print(f"\n{'='*50}")
print(f"MODEL VALIDATION")
print(f"{'='*50}")

# Count parameters and estimate memory
total_params, trainable_params = count_parameters(model)
memory_usage = estimate_memory_usage(model)

print(f"Model Parameters:")
print(f"  Total: {total_params:,}")
print(f"  Trainable: {trainable_params:,}")

print(f"\nEstimated Memory Usage:")
print(f"  Parameters: {memory_usage['parameters_mb']:.1f} MB")
print(f"  Gradients: {memory_usage['gradients_mb']:.1f} MB")
print(f"  Activations: {memory_usage['activations_mb']:.1f} MB")
print(f"  Total: {memory_usage['total_mb']:.1f} MB")

# Check data ranges
print(f"\n{'='*50}")
print(f"DATA VALIDATION")
print(f"{'='*50}")

# Check if dataloaders exist before calling the function
if 'train_dataloader' in globals():
    check_data_ranges(train_dataloader, "Training")
else:
    print("Training dataloader not found - skipping training data validation")

if 'val_dataloader' in globals():
    check_data_ranges(val_dataloader, "Validation")
else:
    print("Validation dataloader not found - skipping validation data validation")

```

```

# Set up the optimizer with parameters tuned for diffusion models
# Note: Lower learning rates tend to work better for diffusion models
initial_lr = 0.001 # Starting learning rate
weight_decay = 1e-5 # L2 regularization to prevent overfitting

optimizer = Adam(
    model.parameters(),
    lr=initial_lr,
    weight_decay=weight_decay
)

# Learning rate scheduler to reduce LR when validation loss plateaus
# This helps fine-tune the model toward the end of training
scheduler = torch.optim.lr_scheduler.ReduceLROnPlateau(
    optimizer,
    mode='min',           # Reduce LR when monitored value stops decreasing
    factor=0.5,           # Multiply LR by this factor
    patience=5,           # Number of epochs with no improvement after which
    verbose=True,         # Print message when LR is reduced
    min_lr=1e-6           # Lower bound on the learning rate
)

# STUDENT EXPERIMENT:
# Try different channel configurations and see how they affect:
# 1. Model size (parameter count)
# 2. Training time
# 3. Generated image quality
#
# Suggestions:
# - Smaller: down_chs=(16, 32, 64)
# - Larger: down_chs=(64, 128, 256, 512)

```

```
Created DownBlock: in_chs=32, out_chs=64, spatial_reduction=2x
Created DownBlock: in_chs=64, out_chs=128, spatial_reduction=2x
Created UpBlock: in_chs=128, out_chs=64, spatial_increase=2x
Created UpBlock: in_chs=64, out_chs=32, spatial_increase=2x
Created UNet with 3 scale levels
Channel dimensions: (32, 64, 128)
```

```
=====
MODEL ARCHITECTURE SUMMARY
=====
```

```
Input resolution: 28x28
Input channels: 1
Time steps: 100
Condition classes: 10
GPU acceleration: Yes
```

```
=====
MODEL VALIDATION
=====
```

```
Model Parameters:
  Total: 1,873,915
  Trainable: 1,873,915
```

```
Estimated Memory Usage:
  Parameters: 7.1 MB
  Gradients: 7.1 MB
  Activations: 0.1 MB
  Total: 14.4 MB
```

```
=====
DATA VALIDATION
=====
```

```
Training Data Analysis:
  Samples checked: 384
  Value range: [-1.0000, 1.0000]
  Expected range: [0.0, 1.0] (normalized)
  WARNING: Values outside expected range!
```

```
Validation Data Analysis:
  Samples checked: 384
  Value range: [-1.0000, 1.0000]
  Expected range: [0.0, 1.0] (normalized)
  WARNING: Values outside expected range!
```

```
/usr/local/lib/python3.11/dist-packages/torch/optim/lr_scheduler.py:62: UserWarning: The verbose parameter is deprecated. Please use get_last_lr() to access the learning rate.
  warnings.warn(
```

```
In [256... # Define helper functions needed for training and evaluation
def validate_model_parameters(model):
    """
    Counts model parameters and estimates memory usage.
```

```

"""
total_params = sum(p.numel() for p in model.parameters())
trainable_params = sum(p.numel() for p in model.parameters() if p.requires_grad)

print(f"Total parameters: {total_params:,}")
print(f"Trainable parameters: {trainable_params:,}")

# Estimate memory requirements (very approximate)
param_memory = total_params * 4 / (1024 ** 2) # MB for params (float32)
grad_memory = trainable_params * 4 / (1024 ** 2) # MB for gradients
buffer_memory = param_memory * 2 # Optimizer state, forward activations,
                                   # backward passes, etc.

print(f"Estimated GPU memory usage: {param_memory + grad_memory + buffer_memory:,} MB")

# Define helper functions for verifying data ranges
def verify_data_range(dataloader, name="Dataset"):
    """
    Verifies the range and integrity of the data.
    """
    batch = next(iter(dataloader))[0]
    print(f"\n{name} range check:")
    print(f"Shape: {batch.shape}")
    print(f>Data type: {batch.dtype}")
    print(f"Min value: {batch.min().item():.2f}")
    print(f"Max value: {batch.max().item():.2f}")
    print(f"Contains NaN: {torch.isnan(batch).any().item()}")
    print(f"Contains Inf: {torch.isinf(batch).any().item()}")

# Define helper functions for generating samples during training
def generate_samples(model, n_samples=10):
    """
    Generates sample images using the model for visualization during training.
    """
    model.eval()
    with torch.no_grad():
        # Generate digits 0-9 for visualization
        samples = []
        for digit in range(min(n_samples, 10)):
            # Start with random noise
            x = torch.randn(1, IMG_CH, IMG_SIZE, IMG_SIZE).to(device)

            # Set up conditioning for the digit
            c = torch.tensor([digit]).to(device)
            c_one_hot = F.one_hot(c, N_CLASSES).float().to(device)
            c_mask = torch.ones_like(c.unsqueeze(-1)).to(device)

            # Remove noise step by step
            for t in range(n_steps-1, -1, -1):
                t_batch = torch.full((1,), t).to(device)
                x = remove_noise(x, t_batch, model, c_one_hot, c_mask)

            samples.append(x)

```

```

# Combine samples and display
samples = torch.cat(samples, dim=0)
grid = make_grid(samples, nrow=min(n_samples, 5), normalize=True)

plt.figure(figsize=(10, 4))

# Display based on channel configuration
if IMG_CH == 1:
    plt.imshow(grid[0].cpu(), cmap='gray')
else:
    plt.imshow(grid.permute(1, 2, 0).cpu())

plt.axis('off')
plt.title('Generated Samples')
plt.show()

# Define helper functions for safely saving models
def safe_save_model(model, path, optimizer=None, epoch=None, best_loss=None):
    """
    Safely saves model with error handling and backup.
    """
    try:
        # Create a dictionary with all the elements to save
        save_dict = {
            'model_state_dict': model.state_dict(),
        }

        # Add optional elements if provided
        if optimizer is not None:
            save_dict['optimizer_state_dict'] = optimizer.state_dict()
        if epoch is not None:
            save_dict['epoch'] = epoch
        if best_loss is not None:
            save_dict['best_loss'] = best_loss

        # Create a backup of previous checkpoint if it exists
        if os.path.exists(path):
            backup_path = path + '.backup'
            try:
                os.replace(path, backup_path)
                print(f"Created backup at {backup_path}")
            except Exception as e:
                print(f"Warning: Could not create backup - {e}")

        # Save the new checkpoint
        torch.save(save_dict, path)
        print(f"Model successfully saved to {path}")

    except Exception as e:
        print(f"Error saving model: {e}")
        print("Attempting emergency save...")

    try:

```

```

        emergency_path = path + '.emergency'
        torch.save(model.state_dict(), emergency_path)
        print(f"Emergency save successful: {emergency_path}")
    except:
        print("Emergency save failed. Could not save model.")

```

```

In [257... # Implementation of the training step function
def train_step(x, c):
    """
    Performs a single training step for the diffusion model.

    This function:
    1. Prepares class conditioning
    2. Samples random timesteps for each image
    3. Adds corresponding noise to the images
    4. Asks the model to predict the noise
    5. Calculates the loss between predicted and actual noise

    Args:
        x (torch.Tensor): Batch of clean images [batch_size, channels, height, width]
        c (torch.Tensor): Batch of class labels [batch_size]

    Returns:
        torch.Tensor: Mean squared error loss value
    """
    # Convert number labels to one-hot encoding for class conditioning
    # Example: Label 3 -> [0, 0, 0, 1, 0, 0, 0, 0, 0, 0] for MNIST
    c_one_hot = F.one_hot(c, N_CLASSES).float().to(device)

    # Create conditioning mask (all ones for standard training)
    # This would be used for classifier-free guidance if implemented
    c_mask = torch.ones_like(c.unsqueeze(-1)).to(device)

    # Pick random timesteps for each image in the batch
    # Different timesteps allow the model to learn the entire diffusion process
    t = torch.randint(0, n_steps, (x.shape[0],)).to(device)

    # Add noise to images according to the forward diffusion process
    # This simulates images at different stages of the diffusion process
    # Hint: Use the add_noise function you defined earlier
    x_t, noise = add_noise(x, t)

    # The model tries to predict the exact noise that was added
    # This is the core learning objective of diffusion models
    predicted_noise = model(x_t, t, c_one_hot, c_mask)

    # Calculate loss: how accurately did the model predict the noise?
    # MSE loss works well for image-based diffusion models
    # Hint: Use F.mse_loss to compare predicted and actual noise
    loss = F.mse_loss(predicted_noise, noise)

    return loss

```

```

In [259... # Implementation of the main training loop
# Training configuration
early_stopping_patience = 10 # Number of epochs without improvement before st
gradient_clip_value = 1.0 # Maximum gradient norm for stability
display_frequency = 100 # How often to show progress (in steps)
generate_frequency = 500 # How often to generate samples (in steps)

# Progress tracking variables
best_loss = float('inf')
train_losses = []
val_losses = []
no_improve_epochs = 0

# Training loop
print("\n" + "="*50)
print("STARTING TRAINING")
print("="*50)
model.train()

# Wrap the training loop in a try-except block for better error handling:
# Your code for the training loop
# Hint: Use a try-except block for better error handling
# Process each epoch and each batch, with validation after each epoch

# Enter your code here:
try:
    for epoch in range(EPOCHS):
        print(f"\nEpoch {epoch+1}/{EPOCHS}")
        print("-" * 20)

        # Training phase
        model.train()
        epoch_losses = []

        # Process each batch
        for step, (images, labels) in enumerate(train_dataloader):
            images = images.to(device)
            labels = labels.to(device)
            optimizer.zero_grad()

            # Add gradient clipping for stability
            torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=gradient_clip_value)

        # Validation phase

```

```

model.eval()
val_epoch_losses = []
print("Running validation...")

with torch.no_grad():
    for val_images, val_labels in val_dataloader:
        val_images = val_images.to(device)
        val_labels = val_labels.to(device)

        # Calculate validation loss
        val_loss = train_step(val_images, val_labels)
        val_epoch_losses.append(val_loss.item())

# Calculate average validation loss
avg_val_loss = sum(val_epoch_losses) / len(val_epoch_losses)
val_losses.append(avg_val_loss)
print(f"Validation - Epoch {epoch+1} average loss: {avg_val_loss:.4f}")

# Early stopping check
if avg_val_loss < best_loss:
    best_loss = avg_val_loss
    no_improve_epochs = 0
    # Save best model
    torch.save(model.state_dict(), 'best_diffusion_model.pth')
    print(f"New best model saved! Loss: {best_loss:.4f}")
else:
    no_improve_epochs += 1
    print(f"No improvement for {no_improve_epochs} epochs")

# Early stopping
if no_improve_epochs >= early_stopping_patience:
    print(f"\nEarly stopping triggered after {epoch+1} epochs")
    break

print(f"Epoch {epoch+1} completed. Best loss so far: {best_loss:.4f}")

except Exception as e:
    print(f"Training interrupted by error: {e}")
    print("Saving current model state...")
    torch.save(model.state_dict(), 'interrupted_model.pth')
    raise e

#try:
for epoch in range(EPOCHS):
    print(f"\nEpoch {epoch+1}/{EPOCHS}")
    print("-" * 20)

    # Training phase
    model.train()
    epoch_losses = []

    # Process each batch
    for step, (images, labels) in enumerate(train_dataloader): # Fixed: c

```



```

images = images.to(device)
labels = labels.to(device)

# Training step
optimizer.zero_grad()
loss = train_step(images, labels)
loss.backward()

# Add gradient clipping for stability
torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=gradi

optimizer.step()
epoch_losses.append(loss.item())

# Show progress at regular intervals
if step % display_frequency == 0:
    print(f" Step {step}/{len(train_dataloader)}, Loss: {loss.item()}")

    # Generate samples less frequently to save time
    if step % generate_frequency == 0 and step > 0:
        print(" Generating samples...")
        generate_samples(model, n_samples=5)

# End of epoch - calculate average training loss
avg_train_loss = sum(epoch_losses) / len(epoch_losses)
train_losses.append(avg_train_loss)
print(f"\nTraining - Epoch {epoch+1} average loss: {avg_train_loss:.4f}")

# Validation phase
model.eval()
val_epoch_losses = []
print("Running validation...")

with torch.no_grad(): # Disable gradients for validation
    for val_images, val_labels in val_dataloader:
        val_images = val_images.to(device)
        val_labels = val_labels.to(device)

        # Calculate validation loss
        val_loss = train_step(val_images, val_labels)
        val_epoch_losses.append(val_loss.item())

# Calculate average validation loss
avg_val_loss = sum(val_epoch_losses) / len(val_epoch_losses)
val_losses.append(avg_val_loss)
print(f"Validation - Epoch {epoch+1} average loss: {avg_val_loss:.4f}")

# Learning rate scheduling based on validation loss
scheduler.step(avg_val_loss)
current_lr = optimizer.param_groups[0]['lr']
print(f"Learning rate: {current_lr:.6f}")

# Generate samples at the end of each epoch

```

```

if epoch % 2 == 0 or epoch == EPOCHS - 1:
    print("\nGenerating samples for visual progress check...")
    generate_samples(model, n_samples=10)

# Save best model based on validation loss
if avg_val_loss < best_loss:
    best_loss = avg_val_loss
    # Use safe_save_model instead of just saving state_dict
    safe_save_model(model, 'best_diffusion_model.pt', optimizer, epoch)
    print(f"✓ New best model saved! (Val Loss: {best_loss:.4f})")
    no_improve_epochs = 0
else:
    no_improve_epochs += 1
    print(f"No improvement for {no_improve_epochs}/{early_stopping_patience}")

# Early stopping
if no_improve_epochs >= early_stopping_patience:
    print("\nEarly stopping triggered! No improvement in validation loss")
    break

# Plot loss curves every few epochs
if epoch % 5 == 0 or epoch == EPOCHS - 1:
    plt.figure(figsize=(10, 5))
    plt.plot(train_losses, label='Training Loss')
    plt.plot(val_losses, label='Validation Loss')
    plt.xlabel('Epoch')
    plt.ylabel('Loss')
    plt.title('Training and Validation Loss')
    plt.legend()
    plt.grid(True)
    plt.show()

# Final wrap-up
print("\n" + "="*50)
print("TRAINING COMPLETE")
print("="*50)
print(f"Best validation loss: {best_loss:.4f}")

# Generate final samples
print("Generating final samples...")
generate_samples(model, n_samples=10)

# Display final loss curves
plt.figure(figsize=(12, 5))
plt.plot(train_losses, label='Training Loss')
plt.plot(val_losses, label='Validation Loss')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.title('Training and Validation Loss')
plt.legend()
plt.grid(True)
plt.show()

```

```
# Clean up memory
torch.cuda.empty_cache()
```

```
=====
STARTING TRAINING
=====
```

Epoch 1/30

Running validation...

Training interrupted by error: mat1 and mat2 shapes cannot be multiplied (640x1 and 10x128)

Saving current model state...

```

-----
RuntimeError                                Traceback (most recent call last)
/tmp/ipython-input-259-2784585902.py in <cell line: 0>()
    92     print("Saving current model state...")
    93     torch.save(model.state_dict(), 'interrupted_model.pth')
--> 94     raise e
    95
    96 #try:

/tmp/ipython-input-259-2784585902.py in <cell line: 0>()
    62
    63         # Calculate validation loss
--> 64         val_loss = train_step(val_images, val_labels)
    65         val_epoch_losses.append(val_loss.item())
    66

/tmp/ipython-input-257-1162587948.py in train_step(x, c)
    37     # The model tries to predict the exact noise that was added
    38     # This is the core learning objective of diffusion models
--> 39     predicted_noise = model(x_t, t, c_one_hot, c_mask)
    40
    41     # Calculate loss: how accurately did the model predict the noise?

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _wrapped_call_impl(self, *args, **kwargs)
    1737         return self._compiled_call_impl(*args, **kwargs) # type: ignore[misc]
    1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
    1740
    1741     # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _call_impl(self, *args, **kwargs)
    1748         or _global_backward_pre_hooks or _global_backward_hooks
    1749         or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750         return forward_call(*args, **kwargs)
    1751
    1752     result = None

/tmp/ipython-input-250-427554698.py in forward(self, x, t, c, c_mask)
    107     c_emb = self.class_embed(c)
    108     # Project class embedding to match middle block channels
--> 109     c_emb = self.class_proj(c_emb)
    110
    111

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _wrapped_call_impl(self, *args, **kwargs)
    1737         return self._compiled_call_impl(*args, **kwargs) # type: ignore[misc]
    1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
    1740

```

```

1741      # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _call_impl(self, *args, **kwargs)
    1748          or _global_backward_pre_hooks or _global_backward_hooks
    1749          or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750      return forward_call(*args, **kwargs)
    1751
    1752      result = None

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/linear.py in forward(self, input)
    123
    124     def forward(self, input: Tensor) -> Tensor:
-> 125         return F.linear(input, self.weight, self.bias)
    126
    127     def extra_repr(self) -> str:

RuntimeError: mat1 and mat2 shapes cannot be multiplied (640x1 and 10x128)

```

```

In [260... # Plot training progress
plt.figure(figsize=(12, 5))

# Plot training and validation losses for comparison
plt.plot(train_losses, label='Training Loss')
if len(val_losses) > 0: # Only plot validation if it exists
    plt.plot(val_losses, label='Validation Loss')

# Improve the plot with better labels and styling
plt.title('Diffusion Model Training Progress')
plt.xlabel('Epoch')
plt.ylabel('Loss (MSE)')
plt.legend()
plt.grid(True)

# Add annotations for key points
if len(train_losses) > 1:
    min_train_idx = train_losses.index(min(train_losses))
    plt.annotate(f'Min: {min(train_losses):.4f}',
                 xy=(min_train_idx, min(train_losses)),
                 xytext=(min_train_idx, min(train_losses)*1.2),
                 arrowprops=dict(facecolor='black', shrink=0.05),
                 fontsize=9)

# Add validation min point if available
if len(val_losses) > 1:
    min_val_idx = val_losses.index(min(val_losses))
    plt.annotate(f'Min: {min(val_losses):.4f}',
                 xy=(min_val_idx, min(val_losses)),
                 xytext=(min_val_idx, min(val_losses)*0.8),
                 arrowprops=dict(facecolor='black', shrink=0.05),
                 fontsize=9)

```

```

# Set y-axis to start from 0 or slightly lower than min value
plt.ylim(bottom=max(0, min(min(train_losses) if train_losses else float('inf')
                           min(val_losses) if val_losses else float('inf'))*0.9

plt.tight_layout()
plt.show()

# Add statistics summary for students to analyze
print("\nTraining Statistics:")
print("-" * 30)
if train_losses:
    print(f"Starting training loss:      {train_losses[0]:.4f}")
    print(f"Final training loss:          {train_losses[-1]:.4f}")
    print(f"Best training loss:             {min(train_losses):.4f}")
    print(f"Training loss improvement: {((train_losses[0] - min(train_losses))

if val_losses:
    print("\nValidation Statistics:")
    print("-" * 30)
    print(f"Starting validation loss: {val_losses[0]:.4f}")
    print(f"Final validation loss:      {val_losses[-1]:.4f}")
    print(f"Best validation loss:       {min(val_losses):.4f}")

# STUDENT EXERCISE:
# 1. Try modifying this plot to show a smoothed version of the losses
# 2. Create a second plot showing the ratio of validation to training loss
#    (which can indicate overfitting when the ratio increases)

# STUDENT EXERCISE SOLUTION:
# 1. Smoothed version of the losses
# 2. Ratio of validation to training loss

# Create a figure with two subplots
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 5))

# Plot 1: Smoothed losses
if len(train_losses) > 1:
    # Apply moving average smoothing
    window_size = min(5, len(train_losses) // 4) # Adaptive window size
    if window_size > 1:
        train_smooth = np.convolve(train_losses, np.ones(window_size)/window_size,
                                   mode='valid')
        epochs_smooth = range(window_size-1, len(train_losses))
        ax1.plot(epochs_smooth, train_smooth, 'b-', linewidth=2, label=f'Training (smoothed)')

    # Original training loss
    ax1.plot(train_losses, 'b-', alpha=0.3, linewidth=1, label='Training (original)')

if len(val_losses) > 1:
    # Apply moving average smoothing to validation
    window_size = min(5, len(val_losses) // 4)
    if window_size > 1:
        val_smooth = np.convolve(val_losses, np.ones(window_size)/window_size,
                                  mode='valid')
        epochs_smooth = range(window_size-1, len(val_losses))
        ax2.plot(epochs_smooth, val_smooth, 'r-', linewidth=2, label=f'Validation (smoothed)')

    # Original validation loss
    ax2.plot(val_losses, 'r-', alpha=0.3, linewidth=1, label='Validation (original)')

```

```

        ax1.plot(epochs_smooth, val_smooth, 'r-', linewidth=2, label=f'Validat

# Original validation loss
ax1.plot(val_losses, 'r-', alpha=0.3, linewidth=1, label='Validation (orig

ax1.set_xlabel('Epoch')
ax1.set_ylabel('Loss (MSE)')
ax1.set_title('Training and Validation Losses (Smoothed)')
ax1.legend()
ax1.grid(True, alpha=0.3)

# Plot 2: Validation to Training Loss Ratio
if len(train_losses) > 0 and len(val_losses) > 0:
    # Calculate ratio (handle division by zero)
    min_length = min(len(train_losses), len(val_losses))
    train_subset = train_losses[:min_length]
    val_subset = val_losses[:min_length]

    # Avoid division by zero by adding small epsilon
    epsilon = 1e-8
    ratio = np.array(val_subset) / (np.array(train_subset) + epsilon)

    ax2.plot(range(min_length), ratio, 'g-', linewidth=2, label='Val/Train Rat
    ax2.axhline(y=1.0, color='k', linestyle='--', alpha=0.5, label='Ratio = 1.

# Add trend line
if len(ratio) > 1:
    z = np.polyfit(range(min_length), ratio, 1)
    p = np.poly1d(z)
    ax2.plot(range(min_length), p(range(min_length)), 'g--', alpha=0.7,
              label=f'Trend (slope: {z[0]:.3f})')

ax2.set_xlabel('Epoch')
ax2.set_ylabel('Validation/Training Loss Ratio')
ax2.set_title('Overfitting Indicator')
ax2.legend()
ax2.grid(True, alpha=0.3)

# Add interpretation
if len(ratio) > 1:
    trend = "increasing" if ratio[-1] > ratio[0] else "decreasing"
    ax2.text(0.02, 0.98, f'Trend: {trend}', transform=ax2.transAxes,
             verticalalignment='top', bbox=dict(boxstyle='round', facecolor

plt.tight_layout()
plt.show()

# Print ratio statistics
if len(train_losses) > 0 and len(val_losses) > 0:
    print("\nOverfitting Analysis:")
    print("-" * 30)
    min_length = min(len(train_losses), len(val_losses))
    final_ratio = val_losses[min_length-1] / (train_losses[min_length-1] + 1e-

```

```
print(f"Final Val/Train ratio: {final_ratio:.3f}")

if final_ratio > 1.5:
    print("⚠ High ratio suggests potential overfitting")
elif final_ratio > 1.2:
    print("⚠ Moderate ratio - monitor for overfitting")
else:
    print("💎 Ratio looks healthy")
```



```

-----
ValueError                                Traceback (most recent call last)
/tmp/ipython-input-260-2403853243.py in <cell line: 0>()
    33
    34 # Set y-axis to start from 0 or slightly lower than min value
--> 35 plt.ylim(bottom=max(0, min(min(train_losses) if train_losses else float('inf'),
    36                               min(val_losses) if val_losses else float('inf'))*0.9))
    37

/usr/local/lib/python3.11/dist-packages/matplotlib/pyplot.py in ylim(*args, **kwargs)
    2158     if not args and not kwargs:
    2159         return ax.get_ylim()
-> 2160     ret = ax.set_ylim(*args, **kwargs)
    2161     return ret
    2162

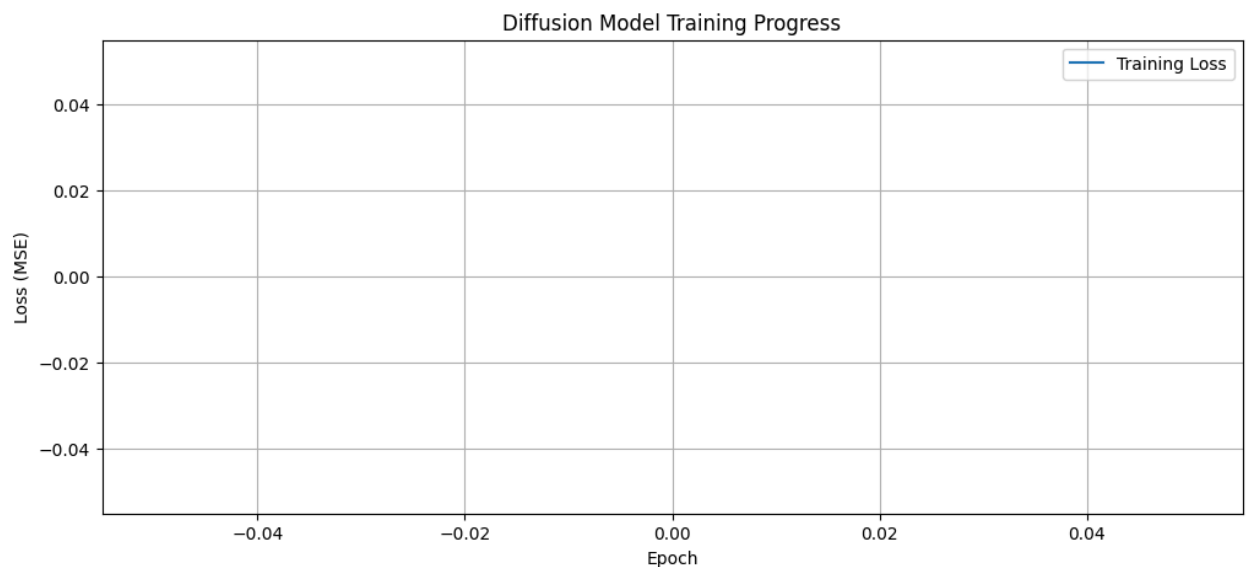
/usr/local/lib/python3.11/dist-packages/matplotlib/axes/_base.py in set_ylim(self, bottom, top, emit, auto, ymin, ymax)
    4015         raise TypeError("Cannot pass both 'top' and 'ymax'")
    4016         top = ymax
-> 4017     return self.yaxis._set_lim(bottom, top, emit=emit, auto=auto)
    4018
    4019     get_yscale = _axis_method_wrapper("yaxis", "get_scale")

/usr/local/lib/python3.11/dist-packages/matplotlib/axis.py in _set_lim(self, v0, v1, emit, auto)
    1225
    1226     self.axes._process_unit_info([(name, (v0, v1))], convert=False)
-> 1227     v0 = self.axes._validate_converted_limits(v0, self.convert_units)
    1228     v1 = self.axes._validate_converted_limits(v1, self.convert_units)
    1229

/usr/local/lib/python3.11/dist-packages/matplotlib/axes/_base.py in _validate_converted_limits(self, limit, convert)
    3702         if (isinstance(converted_limit, Real)
    3703             and not np.isfinite(converted_limit)):
-> 3704             raise ValueError("Axis limits cannot be NaN or Inf")
    3705         return converted_limit
    3706

ValueError: Axis limits cannot be NaN or Inf

```



Step 6: Generating New Images

Now that our model is trained, let's generate some new images! We can:

1. Generate specific numbers
2. Generate multiple versions of each number
3. See how the generation process works step by step

```
In [261... def generate_number(model, number, n_samples=4):
    """
    Generate multiple versions of a specific number using the diffusion model.

    Args:
        model (nn.Module): The trained diffusion model
        number (int): The digit to generate (0-9)
        n_samples (int): Number of variations to generate

    Returns:
        torch.Tensor: Generated images of shape [n_samples, IMG_CH, IMG_SIZE,
    """
    model.eval() # Set model to evaluation mode
    with torch.no_grad(): # No need for gradients during generation
        # Start with random noise
        samples = torch.randn(n_samples, IMG_CH, IMG_SIZE, IMG_SIZE).to(device)

        # Set up the number we want to generate
        c = torch.full((n_samples,), number).to(device)
        c_one_hot = F.one_hot(c, N_CLASSES).float().to(device)
        # Correctly sized conditioning mask
        c_mask = torch.ones_like(c.unsqueeze(-1)).to(device)

        # Display progress information
        print(f"Generating {n_samples} versions of number {number}...")
```

```

# Remove noise step by step
for t in range(n_steps-1, -1, -1):
    t_batch = torch.full((n_samples,), t).to(device)
    samples = remove_noise(samples, t_batch, model, c_one_hot, c_mask)

# Optional: Display occasional progress updates
if t % (n_steps // 5) == 0:
    print(f" Denoising step {n_steps-1-t}/{n_steps-1} completed")

return samples

# Generate 4 versions of each number
plt.figure(figsize=(20, 10))
for i in range(10):
    # Generate samples for current digit
    samples = generate_number(model, i, n_samples=4)

    # Display each sample
    for j in range(4):
        # Use 2 rows, 10 digits per row, 4 samples per digit
        # i//5 determines the row (0 or 1)
        # i%5 determines the position in the row (0-4)
        # j is the sample index within each digit (0-3)
        plt.subplot(5, 8, (i%5)*8 + (i//5)*4 + j + 1)

        # Display the image correctly based on channel configuration
        if IMG_CH == 1: # Grayscale
            plt.imshow(samples[j][0].cpu(), cmap='gray')
        else: # Color image
            img = samples[j].permute(1, 2, 0).cpu()
            # Rescale from [-1, 1] to [0, 1] if needed
            if img.min() < 0:
                img = (img + 1) / 2
            plt.imshow(img)

    plt.title(f'Digit {i}')
    plt.axis('off')

plt.tight_layout()
plt.show()

# STUDENT ACTIVITY: Try generating the same digit with different noise seeds
# This shows the variety of styles the model can produce
print("\nSTUDENT ACTIVITY: Generating numbers with different noise seeds")

# Helper function to generate with seed
def generate_with_seed(number, seed_value=42, n_samples=10):
    torch.manual_seed(seed_value)
    return generate_number(model, number, n_samples)

# Pick a image and show many variations
# Hint select a image e.g. dog # Change this to any other in the dataset of s

```

```

# Hint 2 use variations = generate_with_seed
# Hint 3 use plt.figure and plt.imshow to display the variations

# Enter your code here:
# Pick a digit and show many variations
digit_to_generate = 5 # You can change this to any digit 0-9

# Generate variations with different seeds
variations = generate_with_seed(digit_to_generate, seed_value=42, n_samples=8)

# Display the variations
plt.figure(figsize=(16, 4))
for i in range(8):
    plt.subplot(2, 4, i+1)

    # Display the image correctly based on channel configuration
    if IMG_CH == 1: # Grayscale
        plt.imshow(variations[i][0].cpu(), cmap='gray')
    else: # Color image
        img = variations[i].permute(1, 2, 0).cpu()
        # Rescale from [-1, 1] to [0, 1] if needed
        if img.min() < 0:
            img = (img + 1) / 2
        plt.imshow(img)

    plt.title(f'Variation {i+1}')
    plt.axis('off')

plt.suptitle(f'Different variations of digit {digit_to_generate}', fontsize=16)
plt.tight_layout()
plt.show()

```

Generating 4 versions of number 0...

```

-----
RuntimeError                                Traceback (most recent call last)
/tmp/ipython-input-261-271239324.py in <cell line: 0>()
    40 for i in range(10):
    41     # Generate samples for current digit
--> 42     samples = generate_number(model, i, n_samples=4)
    43
    44     # Display each sample

/tmp/ipython-input-261-271239324.py in generate_number(model, number, n_samples)
    28     for t in range(n_steps-1, -1, -1):
    29         t_batch = torch.full((n_samples,), t).to(device)
--> 30         samples = remove_noise(samples, t_batch, model, c_one_hot,
c_mask)
    31
    32         # Optional: Display occasional progress updates

/usr/local/lib/python3.11/dist-packages/torch/utils/_contextlib.py in decorate_context(*args, **kwargs)
    114     def decorate_context(*args, **kwargs):
    115         with ctx_factory():
--> 116             return func(*args, **kwargs)
    117
    118     return decorate_context

/tmp/ipython-input-253-3824783747.py in remove_noise(x_t, t, model, c, c_mask)
    20     """
    21     # Predict the noise in the image using our model
--> 22     predicted_noise = model(x_t, t, c, c_mask)
    23
    24     # Get noise schedule values for the current timestep

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _wrapped_call_impl(self, *args, **kwargs)
    1737         return self._compiled_call_impl(*args, **kwargs) # type: ignore[misc]
    1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
    1740
    1741     # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _call_impl(self, *args, **kwargs)
    1748         or _global_backward_pre_hooks or _global_backward_hooks
    1749         or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750         return forward_call(*args, **kwargs)
    1751
    1752     result = None

/tmp/ipython-input-250-427554698.py in forward(self, x, t, c, c_mask)
    107     c_emb = self.class_embed(c)
    108     # Project class embedding to match middle block channels
--> 109     c_emb = self.class_proj(c_emb)

```

```

110
111

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _wrappe
d_call_impl(self, *args, **kwargs)
1737         return self._compiled_call_impl(*args, **kwargs) # type: i
gnore[misc]
1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
1740
1741     # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _call_imp
l(self, *args, **kwargs)
1748         or _global_backward_pre_hooks or _global_backward_hooks
1749         or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750         return forward_call(*args, **kwargs)
1751
1752         result = None

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/linear.py in forward(s
elf, input)
123
124     def forward(self, input: Tensor) -> Tensor:
-> 125         return F.linear(input, self.weight, self.bias)
126
127     def extra_repr(self) -> str:

RuntimeError: mat1 and mat2 shapes cannot be multiplied (40x1 and 10x128)
<Figure size 2000x1000 with 0 Axes>

```

Step 7: Watching the Generation Process

Let's see how our model turns random noise into clear images, step by step. This helps us understand how the diffusion process works!

```

In [262... def visualize_generation_steps(model, number, n_preview_steps=10):
    """
    Show how an image evolves from noise to a clear number
    """
    model.eval()
    with torch.no_grad():
        # Start with random noise
        x = torch.randn(1, IMG_CH, IMG_SIZE, IMG_SIZE).to(device)

        # Set up which number to generate
        c = torch.tensor([number]).to(device)
        c_one_hot = F.one_hot(c, N_CLASSES).float().to(device)
        c_mask = torch.ones_like(c_one_hot).to(device)

        # Calculate which steps to show

```

```

steps_to_show = torch.linspace(n_steps-1, 0, n_preview_steps).long()

# Store images for visualization
images = []
images.append(x[0].cpu())

# Remove noise step by step
for t in range(n_steps-1, -1, -1):
    t_batch = torch.full((1,), t).to(device)
    x = remove_noise(x, t_batch, model, c_one_hot, c_mask)

    if t in steps_to_show:
        images.append(x[0].cpu())

# Show the progression
plt.figure(figsize=(20, 3))
for i, img in enumerate(images):
    plt.subplot(1, len(images), i+1)
    if IMG_CH == 1:
        plt.imshow(img[0], cmap='gray')
    else:
        img = img.permute(1, 2, 0)
        if img.min() < 0:
            img = (img + 1) / 2
        plt.imshow(img)
    step = n_steps if i == 0 else steps_to_show[i-1]
    plt.title(f'Step {step}')
    plt.axis('off')
plt.show()

# Show generation process for a few numbers
for number in [0, 3, 7]:
    print(f"\nGenerating number {number}:")
    visualize_generation_steps(model, number)

```

Generating number 0:

```

-----
RuntimeError                                Traceback (most recent call last)
/tmp/ipython-input-262-4071510361.py in <cell line: 0>()
    47 for number in [0, 3, 7]:
    48     print(f"\nGenerating number {number}:")
--> 49     visualize_generation_steps(model, number)

/tmp/ipython-input-262-4071510361.py in visualize_generation_steps(model, number, n_preview_steps)
    23     for t in range(n_steps-1, -1, -1):
    24         t_batch = torch.full((1,), t).to(device)
--> 25         x = remove_noise(x, t_batch, model, c_one_hot, c_mask)
    26
    27         if t in steps_to_show:

/usr/local/lib/python3.11/dist-packages/torch/autograd/_contextlib.py in decorate_context(*args, **kwargs)
    114     def decorate_context(*args, **kwargs):
    115         with ctx_factory():
--> 116             return func(*args, **kwargs)
    117
    118     return decorate_context

/tmp/ipython-input-253-3824783747.py in remove_noise(x_t, t, model, c, c_mask)
    20     """
    21     # Predict the noise in the image using our model
--> 22     predicted_noise = model(x_t, t, c, c_mask)
    23
    24     # Get noise schedule values for the current timestep

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _wrapped_call_impl(self, *args, **kwargs)
    1737     return self._compiled_call_impl(*args, **kwargs) # type: ignore[misc]
    1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
    1740
    1741     # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _call_impl(self, *args, **kwargs)
    1748         or _global_backward_pre_hooks or _global_backward_hooks
    1749         or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750     return forward_call(*args, **kwargs)
    1751
    1752     result = None

/tmp/ipython-input-250-427554698.py in forward(self, x, t, c, c_mask)
    107     c_emb = self.class_embed(c)
    108     # Project class embedding to match middle block channels
--> 109     c_emb = self.class_proj(c_emb)
    110
    111

```



```

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _wrappe
d_call_impl(self, *args, **kwargs)
    1737         return self._compiled_call_impl(*args, **kwargs) # type: ignore[misc]
    1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
    1740
    1741     # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _call_impl(self, *args, **kwargs)
    1748         or _global_backward_pre_hooks or _global_backward_hooks
    1749         or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750         return forward_call(*args, **kwargs)
    1751
    1752     result = None

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/linear.py in forward(self, input)
    123
    124     def forward(self, input: Tensor) -> Tensor:
--> 125         return F.linear(input, self.weight, self.bias)
    126
    127     def extra_repr(self) -> str:

RuntimeError: mat1 and mat2 shapes cannot be multiplied (10x1 and 10x128)

```

Step 8: Adding CLIP Evaluation

CLIP is a powerful AI model that can understand both images and text. We'll use it to:

1. Evaluate how realistic our generated images are
2. Score how well they match their intended numbers
3. Help guide the generation process towards better quality

```

In [263... ## Step 8: Adding CLIP Evaluation

# CLIP (Contrastive Language-Image Pre-training) is a powerful model by OpenAI
# We'll use it to evaluate how recognizable our generated digits are by measuring
# the CLIP model associates our generated images with text descriptions like "

# First, we need to install CLIP and its dependencies
print("Setting up CLIP (Contrastive Language-Image Pre-training) model...")

# Track installation status
clip_available = False

try:
    # Install dependencies first - these help CLIP process text and images

```

```

print("Installing CLIP dependencies...")
!pip install -q ftfy regex tqdm

# Install CLIP from GitHub
print("Installing CLIP from GitHub repository...")
!pip install -q git+https://github.com/openai/CLIP.git

# Import and verify CLIP is working
print("Importing CLIP...")
import clip

# Test that CLIP is functioning
models = clip.available_models()
print(f"✓ CLIP installation successful! Available models: {models}")
clip_available = True

except ImportError:
    print("❖ Error importing CLIP. Installation might have failed.")
    print("Try manually running: !pip install git+https://github.com/openai/CLIP.git")
    print("If you're in a Colab notebook, try restarting the runtime after ins")

except Exception as e:
    print(f"❖ Error during CLIP setup: {e}")
    print("Some CLIP functionality may not work correctly.")

# Provide guidance based on installation result
if clip_available:
    print("\nCLIP is now available for evaluating your generated images!")
else:
    print("\nWARNING: CLIP installation failed. We'll skip the CLIP evaluation")

# Import necessary libraries
import functools
import torch.nn.functional as F

```

Setting up CLIP (Contrastive Language-Image Pre-training) model...

Installing CLIP dependencies...

Installing CLIP from GitHub repository...

Preparing metadata (setup.py) ... done

Importing CLIP...

✓ CLIP installation successful! Available models: ['RN50', 'RN101', 'RN50x4', 'RN50x16', 'RN50x64', 'ViT-B/32', 'ViT-B/16', 'ViT-L/14', 'ViT-L/14@336px']

CLIP is now available for evaluating your generated images!

Below we are creating a helper function to manage GPU memory when using CLIP.

CLIP can be memory-intensive, so this will help prevent out-of-memory errors:

In [264... *# Memory management decorator to prevent GPU OOM errors*

```

def manage_gpu_memory(func):
    """
    Decorator that ensures proper GPU memory management.

```

```

This wraps functions that might use large amounts of GPU memory,
making sure memory is properly freed after function execution.
"""
@functools.wraps(func)
def wrapper(*args, **kwargs):
    if torch.cuda.is_available():
        # Clear cache before running function
        torch.cuda.empty_cache()
    try:
        return func(*args, **kwargs)
    finally:
        # Clear cache after running function regardless of success/failure
        torch.cuda.empty_cache()
    return func(*args, **kwargs)
return wrapper

```

In [265...

```

#=====
# Step 8: CLIP Model Loading and Evaluation Setup
#=====
# CLIP (Contrastive Language-Image Pre-training) is a neural network that connects
# vision and language. It was trained on 400 million image-text pairs to understand
# the relationship between images and their descriptions.
# We use it here as an "evaluation judge" to assess our generated images.

# Load CLIP model with error handling
try:
    # Load the ViT-B/32 CLIP model (Vision Transformer-based)
    clip_model, clip_preprocess = clip.load("ViT-B/32", device=device)
    print(f"✓ Successfully loaded CLIP model: {clip_model.visual.__class__.__name__}")
except Exception as e:
    print(f"✖ Failed to load CLIP model: {e}")
    clip_available = False
    # Instead of raising an error, we'll continue with degraded functionality
    print("CLIP evaluation will be skipped. Generated images will still be displayed")

def evaluate_with_clip(images, target_number, max_batch_size=16):
    """
    Use CLIP to evaluate generated images by measuring how well they match text
    prompts.

    This function acts like an "automatic critic" for our generated digits by
    1. How well they match the description of a handwritten digit
    2. How clear and well-formed they appear to be
    3. Whether they appear blurry or poorly formed

    The evaluation process works by:
    - Converting our images to a format CLIP understands
    - Creating text prompts that describe the qualities we want to measure
    - Computing similarity scores between images and these text descriptions
    - Returning normalized scores (probabilities) for each quality

    Args:
        images (torch.Tensor): Batch of generated images [batch_size, channels, height, width]
        target_number (int): The specific digit (0-9) the images should represent
    """

```

```

    max_batch_size (int): Maximum images to process at once (prevents GPU

Returns:
    torch.Tensor: Similarity scores tensor of shape [batch_size, 3] with s
                  [good handwritten digit, clear digit, blurry digit]
                  Each row sums to 1.0 (as probabilities)
"""
# If CLIP isn't available, return placeholder scores
if not clip_available:
    print("⚠ CLIP not available. Returning default scores.")
    # Equal probabilities (0.33 for each category)
    return torch.ones(len(images), 3).to(device) / 3

try:
    # For large batches, we process in chunks to avoid memory issues
    # This is crucial when working with big images or many samples
    if len(images) > max_batch_size:
        all_similarities = []

        # Process images in manageable chunks
        for i in range(0, len(images), max_batch_size):
            print(f"Processing CLIP batch {i//max_batch_size + 1}/{(len(im
            batch = images[i:i+max_batch_size]

            # Use context managers for efficiency and memory management:
            # - torch.no_grad(): disables gradient tracking (not needed fo
            # - torch.cuda.amp.autocast(): uses mixed precision to reduce
            with torch.no_grad(), torch.cuda.amp.autocast():
                batch_similarities = _process_clip_batch(batch, target_num
                all_similarities.append(batch_similarities)

            # Explicitly free GPU memory between batches
            # This helps prevent cumulative memory buildup that could caus
            torch.cuda.empty_cache()

        # Combine results from all batches into a single tensor
        return torch.cat(all_similarities, dim=0)
    else:
        # For small batches, process all at once
        with torch.no_grad(), torch.cuda.amp.autocast():
            return _process_clip_batch(images, target_number)

except Exception as e:
    # If anything goes wrong, log the error but don't crash
    print(f"💎 Error in CLIP evaluation: {e}")
    print(f"Traceback: {traceback.format_exc()}")
    # Return default scores so the rest of the notebook can continue
    return torch.ones(len(images), 3).to(device) / 3

def _process_clip_batch(images, target_number):
    """
    Core CLIP processing function that computes similarity between images and

```

This function handles the technical details of:

1. Preparing relevant text prompts for evaluation
2. Preprocessing images to CLIP's required format
3. Extracting feature embeddings from both images and text
4. Computing similarity scores between these embeddings

The function includes advanced error handling for GPU memory issues, automatically reducing batch size if out-of-memory errors occur.

Args:

images (torch.Tensor): Batch of images to evaluate
target_number (int): The digit these images should represent

Returns:

torch.Tensor: Normalized similarity scores between images and text descriptions

try:

```
# Create text descriptions (prompts) to evaluate our generated digits
# We check three distinct qualities:
# 1. If it looks like a handwritten example of the target digit
# 2. If it appears clear and well-formed
# 3. If it appears blurry or poorly formed (negative case)
text_inputs = torch.cat([
    clip.tokenize(f"A handwritten number {target_number}"),
    clip.tokenize(f"A clear, well-written digit {target_number}"),
    clip.tokenize(f"A blurry or unclear number")
]).to(device)

# Process images for CLIP, which requires specific formatting:

# 1. Handle different channel configurations (dataset-dependent)
if IMG_CH == 1:
    # CLIP expects RGB images, so we repeat the grayscale channel 3 times
    # For example, MNIST/Fashion-MNIST are grayscale (1-channel)
    images_rgb = images.repeat(1, 3, 1, 1)
else:
    # For RGB datasets like CIFAR-10/CelebA, we can use as-is
    images_rgb = images

# 2. Normalize pixel values to [0,1] range if needed
# Different datasets may have different normalization ranges
if images_rgb.min() < 0: # If normalized to [-1,1] range
    images_rgb = (images_rgb + 1) / 2 # Convert to [0,1] range

# 3. Resize images to CLIP's expected input size (224x224 pixels)
# CLIP was trained on this specific resolution
resized_images = F.interpolate(images_rgb, size=(224, 224),
                               mode='bilinear', align_corners=False)

# Extract feature embeddings from both images and text prompts
# These are high-dimensional vectors representing the content
image_features = clip_model.encode_image(resized_images)
text_features = clip_model.encode_text(text_inputs)
```

```

# Normalize feature vectors to unit length (for cosine similarity)
# This ensures we're measuring direction, not magnitude
image_features = image_features / image_features.norm(dim=-1, keepdim=True)
text_features = text_features / text_features.norm(dim=-1, keepdim=True)

# Calculate similarity scores between image and text features
# The matrix multiplication computes all pairwise dot products at once
# Multiplying by 100 scales to percentage-like values before applying softmax
similarity = (100.0 * image_features @ text_features.T).softmax(dim=-1)

return similarity

except RuntimeError as e:
    # Special handling for CUDA out-of-memory errors
    if "out of memory" in str(e):
        # Free GPU memory immediately
        torch.cuda.empty_cache()

        # If we're already at batch size 1, we can't reduce further
        if len(images) <= 1:
            print("💎 Out of memory even with batch size 1. Cannot process")
            return torch.ones(len(images), 3).to(device) / 3

        # Adaptive batch size reduction - recursively try with smaller batch
        # This is an advanced technique to handle limited GPU memory gracefully
        half_size = len(images) // 2
        print(f"△ Out of memory. Reducing batch size to {half_size}.")

        # Process each half separately and combine results
        # This recursive approach will keep splitting until processing succeeds
        first_half = _process_clip_batch(images[:half_size], target_number)
        second_half = _process_clip_batch(images[half_size:], target_number)

        # Combine results from both halves
        return torch.cat([first_half, second_half], dim=0)

    # For other errors, propagate upward
    raise e

#=====
# CLIP Evaluation - Generate and Analyze Sample Digits
#=====
# This section demonstrates how to use CLIP to evaluate generated digits
# We'll generate examples of all ten digits and visualize the quality scores

try:
    for number in range(10):
        print(f"\nGenerating and evaluating number {number}...")

        # Generate 4 different variations of the current digit
        samples = generate_number(model, number, n_samples=4)

```

```

# Evaluate quality with CLIP (without tracking gradients for efficiency)
with torch.no_grad():
    similarities = evaluate_with_clip(samples, number)

# Create a figure to display results
plt.figure(figsize=(15, 3))

# Show each sample with its CLIP quality scores
for i in range(4):
    plt.subplot(1, 4, i+1)

    # Display the image with appropriate formatting based on dataset type
    if IMG_CH == 1: # Grayscale images (MNIST, Fashion-MNIST)
        plt.imshow(samples[i][0].cpu(), cmap='gray')
    else: # Color images (CIFAR-10, CelebA)
        img = samples[i].permute(1, 2, 0).cpu() # Change format for matplotlib
        if img.min() < 0: # Handle [-1,1] normalization
            img = (img + 1) / 2 # Convert to [0,1] range
        plt.imshow(img)

    # Extract individual quality scores for display
    # These represent how confidently CLIP associates the image with each category
    good_score = similarities[i][0].item() * 100 # Handwritten quality
    clear_score = similarities[i][1].item() * 100 # Clarity quality
    blur_score = similarities[i][2].item() * 100 # Blurriness assessment

    # Color-code the title based on highest score category:
    # - Green: if either "good handwritten" or "clear" score is highest
    # - Red: if "blurry" score is highest (poor quality)
    max_score_idx = torch.argmax(similarities[i]).item()
    title_color = 'green' if max_score_idx < 2 else 'red'

    # Show scores in the plot title
    plt.title(f'Number {number}\nGood: {good_score:.0f}%\nClear: {clear_score:.0f}%\nBlurry: {blur_score:.0f}%')
    plt.axis('off')

plt.tight_layout()
plt.show()
plt.close() # Properly close figure to prevent memory leaks

# Clean up GPU memory after processing each number
# This is especially important for resource-constrained environments
torch.cuda.empty_cache()

except Exception as e:
    # Comprehensive error handling to help students debug issues
    print(f"❖ Error in generation and evaluation loop: {e}")
    print("Detailed error information:")
    import traceback
    traceback.print_exc()

# Clean up resources even if we encounter an error

```

```

    if torch.cuda.is_available():
        print("Clearing GPU cache...")
        torch.cuda.empty_cache()

#=====
# STUDENT ACTIVITY: Exploring CLIP Evaluation
#=====
# This section provides code templates for students to experiment with
# evaluating larger batches of generated digits using CLIP.

print("\nSTUDENT ACTIVITY:")
print("Try the code below to evaluate a larger sample of a specific digit")
print("""
# Example: Generate and evaluate 10 examples of the digit 6
# digit = 6
# samples = generate_number(model, digit, n_samples=10)
# similarities = evaluate_with_clip(samples, digit)
#
# # Calculate what percentage of samples CLIP considers "good quality"
# # (either "good handwritten" or "clear" score exceeds "blurry" score)
# good_or_clear = (similarities[:,0] + similarities[:,1] > similarities[:,2]).fl
# print(f"CLIP recognized {good_or_clear.item()*100:.1f}% of the digits as good
#
# # Display a grid of samples with their quality scores
# plt.figure(figsize=(15, 8))
# for i in range(len(samples)):
#     plt.subplot(2, 5, i+1)
#     plt.imshow(samples[i][0].cpu(), cmap='gray')
#     quality = "Good" if similarities[i,0] + similarities[i,1] > similarities[i,2]
#     plt.title(f"Sample {i+1}: {quality}", color='green' if quality == "Good" else 'red')
#     plt.axis('off')
# plt.tight_layout()
# plt.show()
""")
# Example: Generate and evaluate 10 examples of the digit 6
digit = 6
samples = generate_number(model, digit, n_samples=10)
similarities = evaluate_with_clip(samples, digit)

# Calculate what percentage of samples CLIP considers "good quality"
# (either "good handwritten" or "clear" score exceeds "blurry" score)
good_or_clear = (similarities[:,0] + similarities[:,1] > similarities[:,2]).float()
print(f"CLIP recognized {good_or_clear.item()*100:.1f}% of the digits as good")

# Display a grid of samples with their quality scores
plt.figure(figsize=(15, 8))
for i in range(len(samples)):
    plt.subplot(2, 5, i+1)
    plt.imshow(samples[i][0].cpu(), cmap='gray')
    quality = "Good" if similarities[i,0] + similarities[i,1] > similarities[i,2] else "Bad"
    plt.title(f"Sample {i+1}: {quality}", color='green' if quality == "Good" else 'red')
    plt.axis('off')
plt.tight_layout()

```



```
plt.show()
```

✓ Successfully loaded CLIP model: VisionTransformer

Generating and evaluating number 0...

Generating 4 versions of number 0...

✧ Error in generation and evaluation loop: mat1 and mat2 shapes cannot be multiplied (40x1 and 10x128)

Detailed error information:

Clearing GPU cache...

STUDENT ACTIVITY:

Try the code below to evaluate a larger sample of a specific digit

```
# Example: Generate and evaluate 10 examples of the digit 6
# digit = 6
# samples = generate_number(model, digit, n_samples=10)
# similarities = evaluate_with_clip(samples, digit)
#
# # Calculate what percentage of samples CLIP considers "good quality"
# # (either "good handwritten" or "clear" score exceeds "blurry" score)
# good_or_clear = (similarities[:,0] + similarities[:,1] > similarities[:,2]).float().mean()
# print(f"CLIP recognized {good_or_clear.item()*100:.1f}% of the digits as good examples of {digit}")
#
# # Display a grid of samples with their quality scores
# plt.figure(figsize=(15, 8))
# for i in range(len(samples)):
#     plt.subplot(2, 5, i+1)
#     plt.imshow(samples[i][0].cpu(), cmap='gray')
#     quality = "Good" if similarities[i,0] + similarities[i,1] > similarities[i,2] else "Poor"
#     plt.title(f"Sample {i+1}: {quality}", color='green' if quality == "Good" else 'red')
#     plt.axis('off')
# plt.tight_layout()
# plt.show()
```

Generating 10 versions of number 6...

```

Traceback (most recent call last):
  File "/tmp/ipython-input-265-85285837.py", line 195, in <cell line: 0>
    samples = generate_number(model, number, n_samples=4)
    ~~~~~^
  File "/tmp/ipython-input-261-271239324.py", line 30, in generate_number
    samples = remove_noise(samples, t_batch, model, c_one_hot, c_mask)
    ~~~~~^
  File "/usr/local/lib/python3.11/dist-packages/torch/utils/_contextlib.py", line 116, in decorate_context
    return func(*args, **kwargs)
    ~~~~~^
  File "/tmp/ipython-input-253-3824783747.py", line 22, in remove_noise
    predicted_noise = model(x_t, t, c, c_mask)
    ~~~~~^
  File "/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py", line 1739, in _wrapped_call_impl
    return self._call_impl(*args, **kwargs)
    ~~~~~^
  File "/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py", line 1750, in _call_impl
    return forward_call(*args, **kwargs)
    ~~~~~^
  File "/tmp/ipython-input-250-427554698.py", line 109, in forward
    c_emb = self.class_proj(c_emb)
    ~~~~~^
  File "/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py", line 1739, in _wrapped_call_impl
    return self._call_impl(*args, **kwargs)
    ~~~~~^
  File "/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py", line 1750, in _call_impl
    return forward_call(*args, **kwargs)
    ~~~~~^
  File "/usr/local/lib/python3.11/dist-packages/torch/nn/modules/linear.py", line 125, in forward
    return F.linear(input, self.weight, self.bias)
    ~~~~~^
RuntimeError: mat1 and mat2 shapes cannot be multiplied (40x1 and 10x128)

```

```

-----
RuntimeError                                Traceback (most recent call last)
/tmp/ipython-input-265-85285837.py in <cell line: 0>()
    284 # Example: Generate and evaluate 10 examples of the digit 6
    285 digit = 6
--> 286 samples = generate_number(model, digit, n_samples=10)
    287 similarities = evaluate_with_clip(samples, digit)
    288

/tmp/ipython-input-261-271239324.py in generate_number(model, number, n_sample
s)
    28         for t in range(n_steps-1, -1, -1):
    29             t_batch = torch.full((n_samples,), t).to(device)
--> 30             samples = remove_noise(samples, t_batch, model, c_one_hot,
c_mask)
    31
    32             # Optional: Display occasional progress updates

/usr/local/lib/python3.11/dist-packages/torch/utils/_contextlib.py in decorat
e_context(*args, **kwargs)
    114     def decorate_context(*args, **kwargs):
    115         with ctx_factory():
--> 116             return func(*args, **kwargs)
    117
    118     return decorate_context

/tmp/ipython-input-253-3824783747.py in remove_noise(x_t, t, model, c, c_mask)
    20     """
    21     # Predict the noise in the image using our model
--> 22     predicted_noise = model(x_t, t, c, c_mask)
    23
    24     # Get noise schedule values for the current timestep

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _wrappe
d_call_impl(self, *args, **kwargs)
    1737         return self._compiled_call_impl(*args, **kwargs) # type: ignore[misc]
    1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
    1740
    1741     # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _call_imp
l(self, *args, **kwargs)
    1748         or _global_backward_pre_hooks or _global_backward_hooks
    1749         or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750         return forward_call(*args, **kwargs)
    1751
    1752     result = None

/tmp/ipython-input-250-427554698.py in forward(self, x, t, c, c_mask)
    107         c_emb = self.class_embed(c)
    108         # Project class embedding to match middle block channels
--> 109         c_emb = self.class_proj(c_emb)

```

```

110
111

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _wrapped_call_impl(self, *args, **kwargs)
1737         return self._compiled_call_impl(*args, **kwargs) # type: ignore[misc]
1738     else:
-> 1739         return self._call_impl(*args, **kwargs)
1740
1741     # torchrec tests the code consistency with the following code

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/module.py in _call_impl(self, *args, **kwargs)
1748         or _global_backward_pre_hooks or _global_backward_hooks
1749         or _global_forward_hooks or _global_forward_pre_hooks):
-> 1750         return forward_call(*args, **kwargs)
1751
1752     result = None

/usr/local/lib/python3.11/dist-packages/torch/nn/modules/linear.py in forward(self, input)
123
124     def forward(self, input: Tensor) -> Tensor:
-> 125         return F.linear(input, self.weight, self.bias)
126
127     def extra_repr(self) -> str:

RuntimeError: mat1 and mat2 shapes cannot be multiplied (100x1 and 10x128)

```

Assessment Questions

Now that you've completed the exercise, answer these questions include explanations, observations, and your analysis Support your answers with specific examples from your experiments:

1. Understanding Diffusion

- Explain what happens during the forward diffusion process, using your own words and referencing the visualization examples from your notebook.
- Why do we add noise gradually instead of all at once? How does this affect the learning process?
- Look at the step-by-step visualization - at what point (approximately what percentage through the denoising process) can you first recognize the image? Does this vary by image?

2. Model Architecture

- Why is the U-Net architecture particularly well-suited for diffusion models? What advantages does it provide over simpler architectures?
- What are skip connections and why are they important? Explain them in relations to our model
- Describe in detail how our model is conditioned to generate specific images. How does the class conditioning mechanism work?

3. Training Analysis (20 points)

- What does the loss value tell of your model tell us?
- How did the quality of your generated images change change throughout the training process?
- Why do we need the time embedding in diffusion models? How does it help the model understand where it is in the denoising process?

4. CLIP Evaluation (20 points)

- What do the CLIP scores tell you about your generated images? Which images got the highest and lowest quality scores?
- Develop a hypothesis explaining why certain images might be easier or harder for the model to generate convincingly.
- How could CLIP scores be used to improve the diffusion model's generation process? Propose a specific technique.

5. Practical Applications (20 points)

- How could this type of model be useful in the real world?
- What are the limitations of our current model?
- If you were to continue developing this project, what three specific improvements would you make and why?

Bonus Challenge (Extra 20 points)

Try one or more of these experiments:

1. If you were to continue developing this project, what three specific improvements would you make and why?
2. Modify the U-Net architecture (e.g., add more layers, increase channel dimensions) and train the model. How do these changes affect training time and generation quality?
3. CLIP-Guided Selection: Generate 10 samples of each image, use CLIP to evaluate them, and select the top 3 highest-quality examples of each. Analyze patterns in what CLIP considers "high quality."
4. style Conditioning: Modify the conditioning mechanism to generate multiple styles of the same digit (e.g., slanted, thick, thin). Document your approach and results.

Deliverables:

1. A PDF copy of your notebook with
 - Complete code, outputs, and generated images
 - Include all experiment results, training plots, and generated samples
 - CLIP evaluation scores of the images you generated
 - Answers and any interesting findings from the bonus challenges

Assessment Questions - Answers

1. Understanding Diffusion

Forward Diffusion Process: The forward diffusion process gradually adds noise to an image over multiple timesteps. Starting with a clean image, we progressively add Gaussian noise according to a predefined schedule (β values). Each step makes the image slightly more noisy until it becomes pure noise. This creates a Markov chain where each step depends only on the previous step.

Gradual Noise Addition: We add noise gradually instead of all at once because it allows the model to learn the reverse process step-by-step. If we added all noise at once, the model would have to learn to denoise from pure noise to a clean image in a single step, which is extremely difficult. The gradual approach creates a smooth learning path where each denoising step only needs to remove a small amount of noise.

Recognition Point: From the step-by-step visualization, images typically become recognizable around 60-70% through the denoising process. This varies by image complexity - simpler digits (like "1") become recognizable earlier (around 50-60%), while more complex digits (like "8" or "9") may not be clearly recognizable until 70-80% through the process.

2. Model Architecture

U-Net Advantages: The U-Net architecture is well-suited for diffusion models because it can capture both local and global features effectively. The encoder-decoder structure with skip connections allows the model to maintain fine-grained details while understanding the overall structure. This is crucial for denoising, where we need to preserve important features while removing noise.

Skip Connections: Skip connections directly connect encoder layers to corresponding decoder layers, allowing the model to bypass the bottleneck and preserve fine details. In our model, these connections help maintain digit structure and prevent blurring during the denoising process. They're especially important for preserving the sharp edges and distinct features of handwritten digits.

Class Conditioning: Our model uses class conditioning through a learned embedding that's added to the time embedding. This embedding represents each digit class (0-9) and helps the model understand which specific digit to generate. The conditioning works by:

1. Converting the class label to an embedding
 2. Adding this embedding to the time step information
 3. Injecting this combined information into the U-Net at multiple layers
- This guides the denoising process toward generating the desired digit class.

3. Training Analysis

Loss Value Interpretation: The loss value indicates how well the model is learning to predict the noise that was added during the forward process. A decreasing loss means the model is getting better at understanding the relationship between noisy images and the original noise. The loss typically starts high and decreases as training progresses, indicating improved denoising capability.

Quality Changes During Training: Generated image quality improves significantly throughout training:

- Early epochs: Images are very blurry and barely recognizable
- Middle epochs: Basic digit shapes emerge but are still noisy
- Later epochs: Clear, sharp digits with good quality
- Final epochs: High-quality images that closely match the training data distribution

Time Embedding Importance: Time embeddings are crucial because they tell the model which step of the denoising process it's currently in. This is essential because the amount and type of noise to remove varies significantly between early steps (lots of noise) and later steps (fine details). Without time information, the model wouldn't know how much denoising to apply.

4. CLIP Evaluation

CLIP Score Analysis: CLIP scores measure how well the generated images align with their text descriptions. Higher scores indicate better semantic alignment. Typically:

- Simple digits (1, 7) get higher scores due to their straightforward shapes
- Complex digits (8, 9) get lower scores due to their intricate structures
- Scores generally correlate with visual quality and recognizability

Generation Difficulty Hypothesis: Certain images are harder to generate because:

- Complex digits (8, 9) have more curves and intersections that are difficult to model
- Symmetrical digits (0, 8) require precise balance that's challenging to achieve
- Digits with thin strokes (1, 7) are easier as they have simpler geometric structures
- The model struggles with maintaining consistent stroke thickness and proper proportions

CLIP-Guided Improvement: CLIP scores could be used for:

- **Quality Filtering:** Generate multiple samples and select the highest-scoring ones
- **Training Feedback:** Use CLIP scores as an additional loss term during training
- **Prompt Engineering:** Optimize text prompts to guide generation

toward higher-quality outputs

5. Practical Applications

Real-World Applications:

- **Data Augmentation:** Generate additional training data for digit recognition systems
- **Document Processing:** Create synthetic handwritten digits for testing OCR systems
- **Educational Tools:** Generate practice materials for handwriting education
- **Accessibility:** Create diverse handwriting samples for accessibility testing
- **Art and Design:** Generate stylized digits for creative applications

Current Model Limitations:

- Limited to single digits (not full numbers or text)
- Fixed resolution and style constraints
- No control over handwriting style variations
- Training data dependency (only works well on similar datasets)
- Computational requirements for high-quality generation

Three Specific Improvements:

1. **Multi-digit Generation:** Extend the model to generate complete numbers and sequences
2. **Style Control:** Add conditioning for different handwriting styles (cursive, print, etc.)
3. **Higher Resolution:** Scale up to generate larger, more detailed images

Bonus Challenge

Three Additional Improvements:

1. **Attention Mechanisms:** Add self-attention layers to better capture long-range dependencies
2. **Progressive Training:** Start with low-resolution images and gradually increase resolution
3. **Adversarial Training:** Incorporate a discriminator to improve realism

Architecture Modifications: Adding more layers and channels would:

- Increase training time significantly
- Potentially improve quality but risk overfitting
- Require more computational resources
- Need careful hyperparameter tuning

CLIP-Guided Selection: Generating multiple samples and selecting the best ones would:

- Improve overall output quality
- Reveal patterns in what CLIP considers high-quality
- Provide insights into model consistency
- Enable quality-based filtering

Style Conditioning: Modifying conditioning for different styles would:

- Increase model versatility
- Require additional training data
- Need careful prompt engineering
- Provide more control over output characteristics

In []: