

**Exercise 1.**

- (a) We have  $P(\mathbf{x}) = P(x_1)P(x_2|x_1)P(x_3|x_2)P(x_4|x_3)$ .
- (b) Note that by the definition of the  $x_i$ , and letting  $\omega_t \sim \mathcal{N}(0, \sigma^2)$ , we have

$$\begin{aligned} X_1 &= \omega_1 \\ X_2 &= X_1 + \omega_2 \\ X_3 &= X_2 + \omega_3 \\ X_4 &= X_3 + \omega_4 \end{aligned}$$

from which we see that the  $X_i$  are linear combinations of the  $\omega_i$ :

$$X_i = \sum_{1 \leq j \leq i} \omega_j$$

We can represent the  $X_i$  compactly as

$$\begin{aligned} X &= \begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{pmatrix} \\ &= \begin{pmatrix} \omega_1 \\ \omega_1 + \omega_2 \\ \omega_1 + \omega_2 + \omega_3 \\ \omega_1 + \omega_2 + \omega_3 + \omega_4 \end{pmatrix} \\ &= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \omega_1 \\ \omega_2 \\ \omega_3 \\ \omega_4 \end{pmatrix} \\ &=: A\boldsymbol{\omega} \end{aligned}$$

As a result, the covariance matrix is

$$\begin{aligned} A \cdot \text{cov}(\boldsymbol{\omega}) \cdot A^\top &= A \cdot \sigma^2 I_4 \cdot A^\top \\ &= \sigma^2 A \cdot A^\top \\ &= \sigma^2 \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \\ &= \sigma^2 \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 2 & 2 \\ 1 & 2 & 3 & 3 \\ 1 & 2 & 3 & 4 \end{pmatrix} \end{aligned}$$

- (c) Using standard techniques for matrix inversion, we compute the following inverse:

$$A^{-1} = \frac{1}{\sigma^2} \begin{pmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 1 \end{pmatrix}$$

- (d) We observe that  $A_{ij}^{-1} = 0$  iff  $i \neq j$  and  $x_i$  and  $x_j$  are not adjacent in (the undirected version of) graph.

```
Untitled-1

import torch
import random
import numpy as np
import torch.nn as nn
from ResUNet import ConditionalUnet

device = torch.device(
    "cuda"
    if torch.cuda.is_available()
    else "mps"
    if torch.mps.is_available()
    else "cpu"
)
print(f"Using device: {device}")

class ConditionalDDPM(nn.Module):
    def __init__(self, dmconfig):
        super().__init__()
        self.dmconfig = dmconfig
        self.loss_fn = nn.MSELoss()
        self.network = ConditionalUnet(
            1, self.dmconfig.num_feat, self.dmconfig.num_classes
        ).to(device)
        (
            self.betas,
            self.alphas,
            self.sqrt_betas,
            self.oneover_sqrt_alphas,
            self.alpha_bars,
        ) = self.precompute_scheduler()

    def precompute_scheduler(self):
        beta_1, beta_T, T = self.dmconfig.beta_1, self.dmconfig.beta_T, self.dmconfig.T
        betas = torch.tensor(
            np.linspace(beta_1, beta_T, T), dtype=torch.float32, device=device
        )

        alphas = 1 - betas
        sqrt_betas = torch.sqrt(betas)
        oneover_sqrt_alphas = 1 / torch.sqrt(alphas)
        alpha_bars = torch.cumprod(alphas, dim=0)

        return (betas, alphas, sqrt_betas, oneover_sqrt_alphas, alpha_bars)

    def scheduler(self, t_s):
        beta_1, beta_T, T = self.dmconfig.beta_1, self.dmconfig.beta_T, self.dmconfig.T
```

```

# ===== #
# YOUR CODE HERE:
#   Inputs:
#     t_s: the input time steps, with shape (B,1).
#   Outputs:
#     one dictionary containing the variance schedule
#     $\beta_t$ along with other potentially useful constants.

idx = t_s - 1 # for indexing
beta_t = self.betas[idx]
sqrt_beta_t = self.sqrt_betas[idx]
alpha_t = self.alphas[idx]
oneover_sqrt_alpha = self.oneover_sqrt_alphas[idx]
alpha_t_bar = self.alpha_bars[idx]
sqrt_alpha_bar = torch.sqrt(alpha_t_bar)
sqrt_oneminus_alpha_bar = torch.sqrt(1 - alpha_t_bar)

# ===== #
return {
    "beta_t": beta_t,
    "sqrt_beta_t": sqrt_beta_t,
    "alpha_t": alpha_t,
    "sqrt_alpha_bar": sqrt_alpha_bar,
    "oneover_sqrt_alpha": oneover_sqrt_alpha,
    "alpha_t_bar": alpha_t_bar,
    "sqrt_oneminus_alpha_bar": sqrt_oneminus_alpha_bar,
}

def forward(self, images, conditions):
    T = self.dmcconfig.T
    noise_loss = None
    # ===== #
    # YOUR CODE HERE:
    #   Complete the training forward process based on the
    #   given training algorithm.
    #   Inputs:
    #     images: real images from the dataset, with size (B,1,28,28).
    #     conditions: condition labels, with size (B). You should
    #                 convert it to one-hot encoded labels with size (B,10)
    #                 before making it as the input of the denoising network.
    #   Outputs:
    #     noise_loss: loss computed by the self.loss_fn function .

    self.network.train()
    B = conditions.shape[0]

    # convert conditions to one-hot vector
    conditions = torch.nn.functional.one_hot(
        conditions, num_classes=self.dmcconfig.num_classes
    ).to(device)

```

```

# set (mask_p * B) conditions to unconditional values
n_uncond = int(self.dmconfig.mask_p * B)
uncond_indices = torch.tensor(random.sample(range(B), n_uncond), device=device)

conditions[uncond_indices] = self.dmconfig.condition_mask_value

# sample times
times = torch.randint(low=1, high=T + 1, size=(B, 1, 1, 1), device=device)

# sample gaussian noise
epsilon = torch.randn(size=images.shape, device=device)

# corrupt images to timesteps
schedule = self.scheduler(times)
sqrt_alpha_bar = schedule["sqrt_alpha_bar"]
sqrt_oneminus_alpha_bar = schedule["sqrt_oneminus_alpha_bar"]

X_t = sqrt_alpha_bar * images + sqrt_oneminus_alpha_bar * epsilon

# normalize times
times = times.type(torch.float) / T

# compute
out = self.network(X_t, times, conditions)

noise_loss = self.loss_fn(out, epsilon)

# ===== #
return noise_loss

def sample(self, conditions, omega):
    T = self.dmconfig.T
    X_t = None
    # ===== #
    # YOUR CODE HERE:
    # Complete the training forward process based on the
    # given sampling algorithm.
    # Inputs:
    #     conditions: condition labels, with size (B). You should
    #                 convert it to one-hot encoded labels with size (B,10)
    #                 before making it as the input of the denoising network.
    #     omega: conditional guidance weight.
    # Outputs:
    #     generated_images

    self.network.eval()

    B = conditions.shape[0]

```

```
# standard normal sample
X_t = torch.randn(
    size=(B, self.dmcconfig.num_channels, *self.dmcconfig.input_dim),
    device=device,
)

# convert conditions to one-hot
conditions = torch.nn.functional.one_hot(
    conditions, self.dmcconfig.num_classes
).to(device)

# prepare unconditional input
unconditional_conditions = (
    torch.ones_like(conditions) * self.dmcconfig.condition_mask_value
)

with torch.no_grad():
    for t in torch.arange(T, 0, -1, device=device):
        if t > 1:
            z = torch.randn(X_t.shape, device=device)
        else:
            z = 0

        # normalize t
        t_norm = torch.ones((B, 1, 1, 1), device=device) * t
        t_norm = t_norm / T

        # get noise
        conditional_noise = self.network(X_t, t_norm, conditions)
        unconditional_noise = self.network(
            X_t, t_norm, unconditional_conditions
        )
        noise_t = (1 + omega) * conditional_noise - omega * unconditional_noise

        # generate next image
        oneover_sqrt_alpha = self.oneover_sqrt_alphas[t - 1]
        oneminus_alpha = 1 - self.alphas[t - 1]
        oneover_sqrt_oneminus_alpha_bar = 1 / torch.sqrt(
            1 - self.alpha_bars[t - 1]
        )
        sqrt_beta = self.sqrt_betas[t - 1]
        X_t = (
            oneover_sqrt_alpha
            * (X_t - oneminus_alpha * oneover_sqrt_oneminus_alpha_bar * noise_t)
            + sqrt_beta * z
        )

# ===== #
generated_images = (X_t * 0.3081 + 0.1307).clamp(
```

```
    0, 1
) # denormalize the output images
return generated_images
```

**Untitled-2**

```
import torch
import numpy as np
import torch.nn as nn

device = torch.device(
    "cuda"
    if torch.cuda.is_available()
    else "mps"
    if torch.mps.is_available()
    else "cpu"
)

class ResConvBlock(nn.Module):
    """
    Basic residual convolutional block
    """

    def __init__(self, in_channels, out_channels):
        super().__init__()
        self.in_channels = in_channels
        self.out_channels = out_channels
        self.conv1 = nn.Sequential(
            nn.Conv2d(in_channels, out_channels, 3, 1, 1),
            nn.BatchNorm2d(out_channels),
            nn.GELU(),
        )
        self.conv2 = nn.Sequential(
            nn.Conv2d(out_channels, out_channels, 3, 1, 1),
            nn.BatchNorm2d(out_channels),
            nn.GELU(),
        )

    def forward(self, x):
        x1 = self.conv1(x)
        x2 = self.conv2(x1)
        if self.in_channels == self.out_channels:
            out = x + x2
        else:
            out = x1 + x2
        return out / np.sqrt(2)

class UnetDown(nn.Module):
    """
    UNet down block (encoding)
    """
```

```
def __init__(self, in_channels, out_channels):
    super(UnetDown, self).__init__()
    layers = [ResConvBlock(in_channels, out_channels), nn.MaxPool2d(2)]
    self.model = nn.Sequential(*layers)

def forward(self, x):
    return self.model(x)

class UnetUp(nn.Module):
    """
    UNet up block (decoding)
    """

    def __init__(self, in_channels, out_channels):
        super(UnetUp, self).__init__()
        layers = [
            nn.ConvTranspose2d(in_channels, out_channels, 2, 2),
            ResConvBlock(out_channels, out_channels),
            ResConvBlock(out_channels, out_channels),
        ]
        self.model = nn.Sequential(*layers)

    def forward(self, x, skip):
        x = torch.cat((x, skip), 1)
        x = self.model(x)
        return x

class EmbedBlock(nn.Module):
    """
    Embedding block to embed time step/condition to embedding space
    """

    def __init__(self, input_dim, emb_dim):
        super(EmbedBlock, self).__init__()
        self.input_dim = input_dim
        layers = [
            nn.Linear(input_dim, emb_dim),
            nn.GELU(),
            nn.Linear(emb_dim, emb_dim),
        ]
        self.layers = nn.Sequential(*layers)

    def forward(self, x):
        # set embedblock untrainable
        # for param in self.layers.parameters():
        #     param.requires_grad = False
        x = x.view(-1, self.input_dim)
        return self.layers(x)
```

```

class FusionBlock(nn.Module):
    """
    Concatenation and fusion block for adding embeddings
    """

    def __init__(self, in_channels, out_channels):
        super(FusionBlock, self).__init__()
        self.layers = nn.Sequential(
            nn.Conv2d(in_channels, out_channels, 1),
            nn.BatchNorm2d(out_channels),
            nn.GELU(),
        )

    def forward(self, x, t, c):
        h, w = x.shape[-2:]
        return self.layers(
            torch.cat([x, t.repeat(1, 1, h, w), c.repeat(1, 1, h, w)], dim=1)
        )

class ConditionalUnet(nn.Module):
    def __init__(self, in_channels, n_feat=128, n_classes=10):
        super(ConditionalUnet, self).__init__()

        self.in_channels = in_channels
        self.n_feat = n_feat
        self.n_classes = n_classes

        # embeddings
        self.timeembed1 = EmbedBlock(1, 2 * n_feat)
        self.timeembed2 = EmbedBlock(1, 1 * n_feat)
        self.conditionembed1 = EmbedBlock(n_classes, 2 * n_feat)
        self.conditionembed2 = EmbedBlock(n_classes, 1 * n_feat)

        # down path for encoding
        self.init_conv = ResConvBlock(in_channels, n_feat)
        self.downblock1 = UnetDown(n_feat, n_feat)
        self.downblock2 = UnetDown(n_feat, 2 * n_feat)
        self.to_vec = nn.Sequential(nn.AvgPool2d(7), nn.GELU())

        # up path for decoding
        self.upblock0 = nn.Sequential(
            nn.ConvTranspose2d(2 * n_feat, 2 * n_feat, 7, 7),
            nn.GroupNorm(8, 2 * n_feat),
            nn.ReLU(),
        )
        self.upblock1 = UnetUp(4 * n_feat, n_feat)
        self.upblock2 = UnetUp(2 * n_feat, n_feat)

```

```

    self.outblock = nn.Sequential(
        nn.Conv2d(2 * n_feat, n_feat, 3, 1, 1),
        nn.GroupNorm(8, n_feat),
        nn.ReLU(),
        nn.Conv2d(n_feat, self.in_channels, 3, 1, 1),
    )

    # fusion blocks
    self.fusion1 = FusionBlock(3 * self.n_feat, self.n_feat)
    self.fusion2 = FusionBlock(6 * self.n_feat, 2 * self.n_feat)
    self.fusion3 = FusionBlock(3 * self.n_feat, self.n_feat)
    self.fusion4 = FusionBlock(3 * self.n_feat, self.n_feat)

def forward(self, x, t, c):
    """
    Inputs:
        x: input images, with size (B,1,28,28)
        t: input time stepss, with size (B,1,1,1)
        c: input conditions (one-hot encoded labels), with size (B,10)
    """
    t, c = t.float(), c.float()

    # time step embedding
    temb1 = self.timeembed1(t).view(-1, self.n_feat * 2, 1, 1)
    temb2 = self.timeembed2(t).view(-1, self.n_feat, 1, 1)

    # condition embedding
    cemb1 = self.conditionembed1(c).view(-1, self.n_feat * 2, 1, 1)
    cemb2 = self.conditionembed2(c).view(-1, self.n_feat, 1, 1)

    # ===== #
    # YOUR CODE HERE:
    # Define the process of computing the output of a
    # this network given the input x, t, and c.
    # The input x, t, c indicate the input image, time step
    # and the condition respectively.
    # A potential format is shown below, feel free to use your own ways to design it.
    down0 = self.init_conv(x)
    down1 = self.downblock1(down0)
    fus1 = self.fusion1(down1, temb2, cemb2)
    down2 = self.downblock2(fus1)
    fus2 = self.fusion2(down2, temb1, cemb1)
    up0 = self.upblock0(self.to_vec(fus2))
    up1 = self.upblock1(up0, fus2)
    fus3 = self.fusion3(up1, temb2, cemb2)
    up2 = self.upblock2(fus3, fus1)
    fus4 = self.fusion4(up2, temb2, cemb2)
    out = self.outblock(torch.cat((fus4, down0), dim=1))
    # ===== #

```

```
return out
```

# Setup

Similar to the previous projects, we will need some code to set up the environment.

First, run this cell that loads the autoreload extension. This allows us to edit .py source files and re-import them into the notebook for a seamless editing and debugging experience.

```
In [1]: %load_ext autoreload  
%autoreload 2
```

## Google Colab Setup

**If you are not using Colab, please just skip this step.**

Run the following cell to mount your Google Drive. Follow the link and sign in to your Google account (the same account you used to store this notebook!).

```
In [2]: # from google.colab import drive  
# drive.mount('/content/drive')
```

Then enter your path of the project (for example,  
/content/drive/MyDrive/ConditionalDDPM)

```
In [3]: # cd /content/drive/MyDrive/Graduate/ECE239_TA/ConditionalDDPM/skeleton
```

We will use GPUs to accelerate our computation in this notebook.

If you are using Colab, go to `Runtime > Change runtime type` and set `Hardware accelerator` to `GPU`. This will reset Colab. **Rerun the top cell to mount your Drive again.**

Run the following to make sure GPUs are enabled:

```
In [4]: # set the device  
import torch  
device = torch.device("cuda" if torch.cuda.is_available() else "mps" if torch  
  
if torch.cuda.is_available():  
    print('Good to go!')  
else:  
    print('Please set GPU!')
```

Good to go!

# Conditional Denoising Diffusion Probabilistic Models

In the lectures, we have learnt about Denoising Diffusion Probabilistic Models (DDPM), as presented in the paper [Denoising Diffusion Probabilistic Models](#). We went through both the training process and test sampling process of DDPM. In this project, you will use conditional DDPM to generate digits based on given conditions. The project is inspired by the paper [Classifier-free Diffusion Guidance](#), which is a following work of DDPM. You are required to use MNIST dataset and the GPU device to complete the project.

## What is a DDPM?

A Denoising Diffusion Probabilistic Model (DDPM) is a type of generative model inspired by the natural diffusion process. In the example of image generation, DDPM works in two main stages:

- Forward Process (Diffusion): It starts with an image sampled from the dataset and gradually adds noise to it step by step, until it becomes completely random noise. In implementation, the forward diffusion process is fixed to a Markov chain that gradually adds Gaussian noise to the data according to a variance schedule  $\beta_1, \dots, \beta_T$ .
- Reverse Process (Denoising): By learning how the noise was added on the image step by step, the model can do the reverse process: start with random noise and step by step, remove this noise to generate an image.

## Training and sampling of DDPM

As proposed in the DDPM paper, the training and sampling process can be concluded in the following steps:

In [5]: `from IPython.display import Image  
Image(filename='pics/DDPM.png', width=800, height=200)`

Out[5]: **Algorithm 1** Training

```

1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
      $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \epsilon, t)\|^2$ 
6: until converged
  
```

**Algorithm 2** Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
  
```

Here we still use the example of image generation.

Algorithm 1 shows the training process of DDPM. Initially, an image  $\mathbf{x}_0$  is sampled from the data distribution  $q(\mathbf{x}_0)$ , i.e. the dataset. Then a time step  $t$  is randomly selected from a uniform distribution across the predefined number of steps  $T$ .

A noise  $\epsilon$  which has the same shape of the image is sampled from a standard normal distribution. According to the equation (4) in the DDPM paper and the new notation:  $q(\mathbf{x}_t | \mathbf{x}_0) = N(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$ ,  $\alpha_t := 1 - \beta_t$  and  $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$ , we can get an intermediate state of the diffusion process:  $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$ . The model takes the  $\mathbf{x}_t$  and  $t$  as inputs, and predict a noise, i.e.  $\epsilon_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon, t)$ . The optimization of the model is done by minimizing the difference between the sampled noise and the model's prediction of noise.

Algorithm 2 shows the sampling process of DDPM, which is the complete procedure for generating an image. This process starts from noise  $x_T$  sampled from a standard normal distribution, and then uses the trained model to iteratively apply denoising for each time step from  $T$  to 1.

## How to control the generation output?

As you may find, the vanilla DDPM can only randomly generate images which are sampled from the learned distribution of the dataset, while in some cases, we are more interested in controlling the generated images. Previous works mainly use an extra trained classifier to guide the diffusion model to generate specific images ([Dhariwal & Nichol \(2021\)](#)). Ho et al. proposed the [Classifier-free Diffusion Guidance](#), which proposes a novel training and sampling method to achieve the conditional generation without extra models besides the diffusion model. Now let's see how it modify the training and sampling pipeline of DDPM.

### Algorithm 1: Conditional training

The training process is shown in the picture below. Some notations are modified in order to follow DDPM.

```
In [6]: from IPython.display import Image
Image(filename='pics/ConDDPM_1.png', width=800, height=240)
```

Out [6]: **Algorithm 1** Joint training a diffusion model with classifier-free guidance

**Require:**  $p_{\text{uncond}}$ : probability of unconditional training

- 1: **repeat**
- 2:    $(\mathbf{x}_0, \mathbf{c}_0) \sim q(\mathbf{x}_0, \mathbf{c}_0)$  ▷ Sample data with conditioning from the dataset
- 3:    $\mathbf{c}_0 \leftarrow \emptyset$  with probability  $p_{\text{uncond}}$  ▷ Randomly discard conditioning to train unconditionally
- 4:    $t \sim \text{Uniform}(\{1, \dots, T\})$  ▷ Sample time steps
- 5:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6:    $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$  ▷ Corrupt data to the sampled time steps
- 7:   Take gradient step on  $\nabla_\theta \|\epsilon_\theta(\mathbf{x}_t, \mathbf{c}_0, t) - \epsilon\|^2$  ▷ Optimization of denoising model
- 8: **until** converged

Compared with the training process of vanilla DDPM, there are several modifications.

- In the training data sampling, besides the image  $\mathbf{x}_0$ , we also sample the condition  $\mathbf{c}_0$  from the dataset (usually the class label).

- There's a probabilistic step to randomly discard the conditions, training the model to generate data both conditionally and unconditionally. Usually we just set the one-hot encoded label as all -1 to discard the conditions.
- When optimizing the model, the condition  $\mathbf{c}_0$  is an extra input.

### Algorithm 2: Conditional sampling

Below is the sampling process of conditional DDPM.

```
In [7]: from IPython.display import Image
Image(filename='pics/ConDDPM_2.png', width=500, height=250)
```

Out[7]: **Algorithm 2** Conditional sampling with classifier-free guidance

**Require:**  $w$ : guidance weight

**Require:**  $\mathbf{c}$ : conditioning information for conditional sampling

- 1:  $\mathbf{x}_T \sim \mathcal{N}(0, I)$
- 2: **for**  $t = T, \dots, 1$  **do**
- 3:      $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
- 4:      $\tilde{\epsilon}_t = (1 + w)\epsilon_\theta(\mathbf{x}_t, \mathbf{c}, t) - w\epsilon_\theta(\mathbf{x}_t, t)$
- 5:      $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \tilde{\epsilon}_t \right) + \sigma_t \mathbf{z}$
- 6: **end for**
- 7: **return**  $\mathbf{x}_0$

Compared with the vanilla DDPM, the key modification is in step 4. Here the algorithm computes a corrected noise estimation,  $\tilde{\epsilon}_t$ , balancing between the conditional prediction  $\epsilon_\theta(\mathbf{x}_t, \mathbf{c}, t)$  and the unconditional prediction  $\epsilon_\theta(\mathbf{x}_t, t)$ . The corrected noise  $\tilde{\epsilon}_t$  is then used to update  $\mathbf{x}_t$  in step 5.

Here we follow the setting of DDPM paper and define  $\sigma_t = \sqrt{\beta_t}$ .

## Conditional generation of digits

Now let's practice it! You will first asked to design a denoising network, and then complete the training and sampling process of this conditional DDPM.

**In this project, by default, we resize all images to a dimension of  $28 \times 28$  and utilize one-hot encoding for class labels. Also, please remember to normalize the time step  $t$  to the range 0-1 before inputting it into the denoising network as it will help the network have a more stable output.**

First we define a configuration class `DMConfig`. This class contains all the settings of the model and experiment that may be useful later.

```
In [8]: from dataclasses import dataclass, field
from typing import List, Tuple
@dataclass
class DMConfig:
    ...
    Define the model and experiment settings here
    ...
    input_dim: Tuple[int, int] = (28, 28) # input image size
    num_channels: int = 1 # input image channels
    condition_mask_value: int = -1 # unconditional condition mask val
    num_classes: int = 10 # number of classes in the dataset
    T: int = 400 # diffusion and denoising steps
    beta_1: float = 1e-4 # variance schedule
    beta_T: float = 2e-2
    mask_p: float = 0.1 # condition drop ratio
    num_feat: int = 64 # basic feature size of the UNet model
    omega: float = 2.0 # conditional guidance weight
    batch_size: int = 256 # training batch size
    epochs: int = 10 # training epochs
    learning_rate: float = 1e-4 # training learning rate
    multi_lr_milestones: List[int] = field(default_factory=lambda: [20]) # learning rate milestones
    multi_lr_gamma: float = 0.1 # learning rate decay ratio
```

Then let's prepare and visualize the dataset:

```
In [9]: from utils import make_dataloader
from torchvision import transforms
import torchvision.utils as vutils
import matplotlib.pyplot as plt

# Define the data preprocessing and configuration
transform = transforms.Compose([
    transforms.ToTensor(),
    transforms.Normalize((0.1307,), (0.3081,))
])
config = DMConfig()

# Create the train and test dataloaders
train_loader = make_dataloader(transform = transform, batch_size = config.batch_size)
test_loader = make_dataloader(transform = transform, batch_size = config.batch_size)

# Visualize the first 100 images
dataiter = iter(train_loader)
images, labels = next(dataiter)
images_subset = images[:100]
grid = vutils.make_grid(images_subset, nrow = 10, normalize = True, padding=2)
plt.figure(figsize=(6, 6))
plt.imshow(grid.numpy().transpose((1, 2, 0)))
plt.axis('off')
plt.show()
```

Using device: cuda  
Using device: cuda



### 1. Denoising network (6 points)

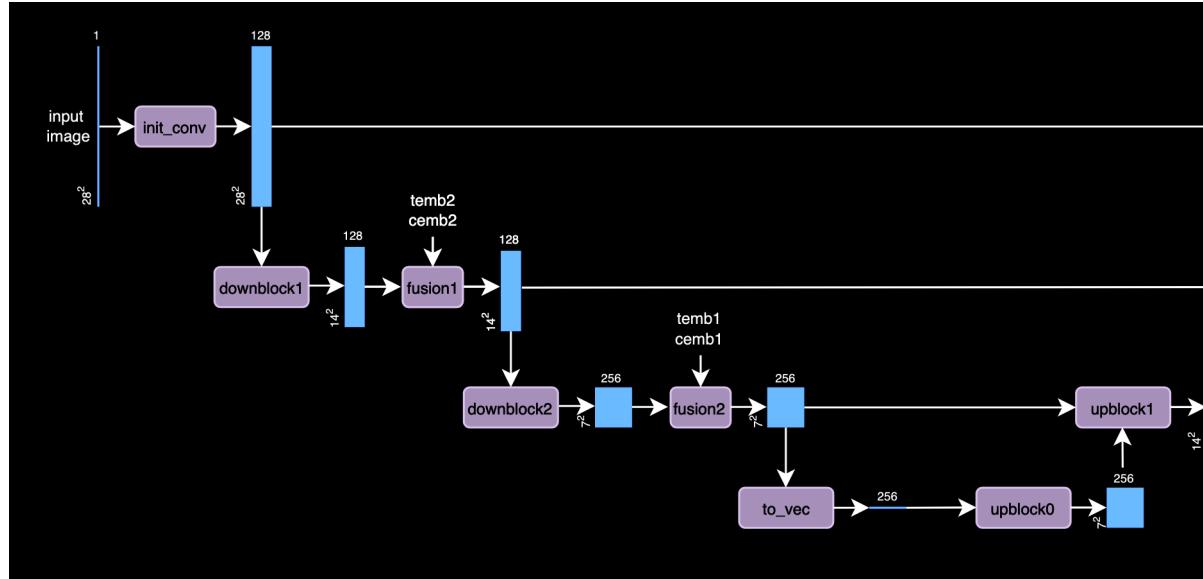
The denoising network is defined in the file `ResUNet.py`. We have already provided some potentially useful layers or blocks, and you will be asked to complete the class `ConditionalUnet`.

Some hints:

- Please consider just using 2 down blocks and 2 up blocks. Using more blocks may improve the performance, while the training and sampling time may increase. Feel free to do some extra experiments in the creative exploring part later.
- An example structure of Conditional UNet is shown in the next cell. Here the initialization argument `n_feat` is set as 128. We provide all the potential useful components in the `__init__` function. The simplest way to construct the network is to complete the `forward` function with these components.
- **MODEL DESIGNING IS AN ART: You do not have to use this given structure. You can add/delete any blocks, or even design your own network from scratch. You are also free to change the way of adding the time step and condition.**

```
In [10]: # Example structure of Conditional UNet
from IPython.core.display import SVG
SVG(filename='./pics/ConUNet.svg')
```

Out[10]:



Now let's check your denoising network using the following code.

```
In [11]: from ResUNet import ConditionalUnet
import torch
model = ConditionalUnet(in_channels = 1, n_feat = 128, n_classes = 10).to(device)
x = torch.randn((256,1,28,28)).to(device)
t = torch.randn((256,1,1,1)).to(device)
c = torch.randn((256,10)).to(device)
x_out = model(x,t,c)
assert x_out.shape == (256,1,28,28)
print('Output shape:', model(x,t,c).shape)
print('Dimension test passed!')
```

Output shape: torch.Size([256, 1, 28, 28])  
Dimension test passed!

## 2. Conditional DDPM

With the correct denoising network, we can then start to build the pipeline of a conditional DDPM. You will be asked to complete the `ConditionalDDPM` class in the file `DDPM.py`.

### 2.1 Variance schedule (4 points)

Let's first prepare the variance schedule  $\beta_t$  along with other potentially useful constants. You are required to complete the `ConditionalDDPM.scheduler` function in `DDPM.py`.

Given the starting and ending variances  $\beta_1$  and  $\beta_T$ , the function should output one dictionary containing the following terms:

`beta_t` : variance of time step  $t_s$ , which is linearly interpolated between  $\beta_1$  and  $\beta_T$ .

`sqrt_beta_t` :  $\sqrt{\beta_t}$

`alpha_t` :  $\alpha_t = 1 - \beta_t$

`oneover_sqrt_alpha` :  $\frac{1}{\sqrt{\alpha_t}}$

`alpha_t_bar` :  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$

`sqrt_alpha_bar` :  $\sqrt{\bar{\alpha}_t}$

`sqrt_oneminus_alpha_bar` :  $\sqrt{1 - \bar{\alpha}_t}$

We set  $\beta_1 = 1e-4$  and  $\beta_T = 2e-2$ . Let's check your solution!

```
In [12]: from DDPM import ConditionalDDPM
torch.set_printoptions(precision=8)
config = DMConfig(beta_1 = 1e-4, beta_T = 2e-2)
ConDDPM = ConditionalDDPM(dmconfig = config)
schedule_dict = ConDDPM.scheduler(t_s = torch.tensor(77)) # We use a specific t_s value here
assert abs(schedule_dict['beta_t'] - 0.003890) <= 5e-6
assert abs(schedule_dict['sqrt_beta_t'] - 0.062374) <= 5e-6
assert abs(schedule_dict['alpha_t'] - 0.996110) <= 5e-6
assert abs(schedule_dict['oneover_sqrt_alpha']) - 1.001951) <= 5e-6
assert abs(schedule_dict['alpha_t_bar'] - 0.857414) <= 5e-6
assert abs(schedule_dict['sqrt_oneminus_alpha_bar'] - 0.377606) <= 5e-6
print('All tests passed!')
```

All tests passed!

## 2.2 Training process (6 points)

Recall the training algorithm we discussed above:

You will need to complete the `ConditionalDDPM.forward` function in the `DDPM.py` file. Then you can use the function `utils.check_forward` to test if it's working properly. The model will be trained for one epoch in this checking process. It should take around 1 min and return one curve showing a decreasing loss trend if your `ConditionalDDPM.forward` function is correct.

```
In [13]: from IPython.display import Image
Image(filename='pics/ConDDPM_1.png', width=800, height=240)
```

---

Out [13]: **Algorithm 1** Joint training a diffusion model with classifier-free guidance

**Require:**  $p_{\text{uncond}}$ : probability of unconditional training

- 1: **repeat**
- 2:    $(\mathbf{x}_0, \mathbf{c}_0) \sim q(\mathbf{x}_0, \mathbf{c}_0)$  ▷ Sample data with conditioning from the dataset
- 3:    $\mathbf{c}_0 \leftarrow \emptyset$  with probability  $p_{\text{uncond}}$  ▷ Randomly discard conditioning to train unconditionally
- 4:    $t \sim \text{Uniform}(\{1, \dots, T\})$  ▷ Sample time steps
- 5:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 6:    $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon$  ▷ Corrupt data to the sampled time steps
- 7:   Take gradient step on  $\nabla_{\theta} \|\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}_0, t) - \epsilon\|^2$  ▷ Optimization of denoising model
- 8: **until** converged

---

In [ ]: `from utils import check_forward  
config = DMConfig()  
model = check_forward(train_loader, config, device)`

### 2.3 Sampling process (6 points)

Now you are required to complete the `ConditionalDDPM.sample` function using the sampling process we mentioned above.

In the following cell, we will use the given `utils.check_sample` function to check the correctness. With the trained model in 2.2, the model should be able to generate some super-rough digits (you may not even see them as digits). The sampling process should take around 30s.

In [15]: `from IPython.display import Image  
Image(filename='pics/ConDDPM_2.png', width=500, height=250)`

---

Out [15]: **Algorithm 2** Conditional sampling with classifier-free guidance

**Require:**  $w$ : guidance weight

**Require:**  $c$ : conditioning information for conditional sampling

- 1:  $\mathbf{x}_T \sim \mathcal{N}(0, I)$
- 2: **for**  $t = T, \dots, 1$  **do**
- 3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = 0$
- 4:    $\tilde{\epsilon}_t = (1 + w)\epsilon_{\theta}(\mathbf{x}_t, \mathbf{c}, t) - w\epsilon_{\theta}(\mathbf{x}_t, t)$
- 5:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\bar{\alpha}_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \tilde{\epsilon}_t \right) + \sigma_t \mathbf{z}$
- 6: **end for**
- 7: **return**  $\mathbf{x}_0$

---

In [ ]: `from utils import check_sample  
config = DMConfig()  
fig = check_sample(model, config, device)`

### 2.4 Full training (8 points)

As you might notice, the images generated are imperfect since the model trained for only one epoch has not yet converged. To improve the performance, we should proceed

with a complete cycle of training and testing. You can utilize a provided `solver` function in this part.

Let's refresh all configurations:

```
In [17]: train_config = DMConfig()
print(train_config)
```

```
DMConfig(input_dim=(28, 28), num_channels=1, condition_mask_value=-1, num_classes=10, T=400, beta_1=0.0001, beta_T=0.02, mask_p=0.1, num_feat=64, omega=2.0, batch_size=256, epochs=10, learning_rate=0.0001, multi_lr_milestones=[20], multi_lr_gamma=0.1)
```

Then we can use function `utils.solver` to train the model. You should also input your own experiment name, e.g. `your_exp_name`. The best-trained model will be saved as `./save/your_exp_name/best_checkpoint.pth`. Furthermore, for each training epoch, one generated image will be stored in the directory `./save/your_exp_name/images` as a validation.

**It will take about 10~20 minutes (10 epochs) if you are using the free-version Google Colab GPU. Typically, realistic digits can be generated after around 2~5 epochs.**

```
In [ ]: from utils import solver
solver(dmconfig = train_config,
       exp_name = 'first_run',
       train_loader = train_loader,
       test_loader = test_loader)
```

**Now please show the image that you believe has the best generation quality in the following cell.**

```
In [19]: # ===== #
# YOUR CODE HERE:
# Among all images generated in the experiment,
# show the image that you believe has the best quality.
# You may use tools like matplotlib, PIL, OpenCV, ...

from PIL import Image
from IPython.display import display

image = Image.open('./save/first_run/images/generate_epoch_9.png')
display(image)

# ===== #
```



### 2.5 Exploring the conditional guidance weight (3 points)

The generated images from the previous training-sampling process is using the default conditional guidance weight  $\omega = 2$ . Now with the best checkpoint, please try at least 3 different  $\omega$  values and visualize the generated images. You can use the provided function `sample_images` to get a combined image each time.

```
In [31]: from utils import sample_images
import matplotlib.pyplot as plt
# ===== #
# YOUR CODE HERE:
# Try at least 3 different conditional guidance weights and visualize it.
# Example of using a different omega value:
#     sample_config = DMConfig(omega = ?)
#     fig = sample_images(config = sample_config, checkpoint_path = path_to_checkpoint)
#     torch.serialization.add_safe_globals([DMConfig])
checkpoint_path = "./save/first_run/best_checkpoint.pth"
omegas = [0.5, 10, 20]
```

```
for omega in omegas:  
    test_config = DMConfig(omega=omega)  
    fig = sample_images(config=test_config, checkpoint_path=checkpoint_path)  
  
    figs, axes = plt.subplots(DMConfig.num_classes, 10, figsize = (6, 6), gr  
    plt.subplots_adjust(wspace=0, hspace=0)  
    axes = axes.flatten()  
    for i, _ in enumerate(fig):  
        ax = axes[i]  
        ax.imshow(fig[i].permute(1,2,0), cmap = 'gray')  
        ax.axis('off')  
    figs.tight_layout()  
  
# ===== #
```







**Inline Question:** Based on your experiment, discuss how the conditional guidance weight affects the quality and diversity of generation. (1 point)

Your answer: Increasing  $\omega$  yields higher quality samples, but decreases sample diversity.

#### 2.6 Customize your own model (5 points)

Now let's experiment by modifying some hyperparameters in the config and customizing your own model. You should at least change one default setting in the config and train a new model. Then visualize the generation image and discuss the effects of your modifications.

**Hint:** Possible changes to the configuration include, but are not limited to, the number of diffusion steps `T`, the unconditional condition drop ratio `mask_p`, the feature size `num_feat`, the beta schedule, etc.

First you should define and print your modified config. Please state all the changes you made to the DMConfig class, i.e. `DMConfig(T=?, num_feat=?, ...)`.

```
In [23]: train_config_new = DMConfig(
    input_dim = (28, 28),
    num_channels = 1,
    condition_mask_value = -1,
    num_classes= 10,
    # T: 400 -> 1000
    T = 1000,
    beta_1 = 0.0001,
    beta_T = 0.02,
    mask_p = 0.1,
    # num_feat: 64 -> 128
    num_feat = 128,
    # omega: 2 -> 3.5
    omega = 3.5,
    batch_size = 256,
    # epochs: 10 -> 12
    epochs = 12,
    learning_rate = 0.0001,
    multi_lr_milestones = [20],
    multi_lr_gamma = 0.1
)
print(train_config_new)
```

```
DMConfig(input_dim=(28, 28), num_channels=1, condition_mask_value=-1, num_classes=10, T=1000, beta_1=0.0001, beta_T=0.02, mask_p=0.1, num_feat=128, omega=3.5, batch_size=256, epochs=12, learning_rate=0.0001, multi_lr_milestones=[20], multi_lr_gamma=0.1)
```

Then similar to 2.4, use `solver` function to complete the training and sampling process.

```
In [24]: from utils import solver
run_name = f'num_feat={train_config_new.num_feat}, epochs={train_config_new.epochs}'
print(run_name)
solver(dmconfig = train_config_new,
       exp_name = f'num_feat={train_config_new.num_feat}, epochs={train_config_new.epochs}',
       train_loader = train_loader,
       test_loader = test_loader)
```

```
num_feat=128, epochs=12, T=1000, omega=3.5
epoch 1/12
```

```
training:  0%|          | 0/235 [00:00<?, ?it/s]
```

```
train: train_noise_loss = 0.1110 test: test_noise_loss = 0.0591
epoch 2/12
```

```
train: train_noise_loss = 0.0524 test: test_noise_loss = 0.0480
epoch 3/12
```

```
train: train_noise_loss = 0.0461 test: test_noise_loss = 0.0456
epoch 4/12
```

```
train: train_noise_loss = 0.0428 test: test_noise_loss = 0.0413
epoch 5/12
```

```
train: train_noise_loss = 0.0404 test: test_noise_loss = 0.0396
epoch 6/12
```

```
train: train_noise_loss = 0.0384 test: test_noise_loss = 0.0418
epoch 7/12
```

```
train: train_noise_loss = 0.0377 test: test_noise_loss = 0.0382
epoch 8/12
```

```
train: train_noise_loss = 0.0364 test: test_noise_loss = 0.0394
epoch 9/12
```

```
train: train_noise_loss = 0.0360 test: test_noise_loss = 0.0378
epoch 10/12
```

```
train: train_noise_loss = 0.0354 test: test_noise_loss = 0.0362
epoch 11/12
```

```
train: train_noise_loss = 0.0355 test: test_noise_loss = 0.0342
epoch 12/12
```

```
train: train_noise_loss = 0.0341 test: test_noise_loss = 0.0357
```

Finally, show one image that you think has the best quality.

```
In [30]: # ===== #
# YOUR CODE HERE:
# Among all images generated in the experiment,
# show the image that you believe has the best generation quality.
# You may use tools like matplotlib, PIL, OpenCV, ...

from PIL import Image
from IPython.display import display

path = './save/num_feat=128, epochs=12, T=1000, omega=3.5/images/generate_еп

image = Image.open(path)
display(image)

# ===== #
```



**Inline Question: Discuss the effects of your modifications after you compare the generation performance under different configurations. (1 point)**

Your answer:

1. Increasing the number of features led to higher sharpness, generating clearer and thinner digits with fewer artifacts.
2. Increasing  $T$  also led to marginally clearer digits.
3. Training for more epochs increased the plausibility of samples.
4. A slightly higher value of omega yielded better samples.