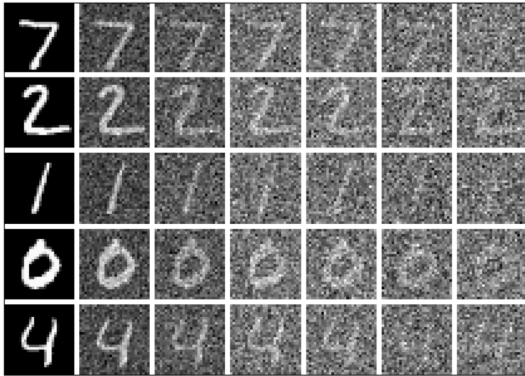


Spring 2025 CS 444

Assignment 5: Diffusion Models

Due date: Wednesday, May 7, 11:59:59 PM

In this assignment, you will train a [Denoising Diffusion Probabilistic Model \(DDPM\)](#) on the MNIST digit dataset.



Credit: This assignment is adapted by Rex Chang from [Fall 2024 CS180 Project 5b](#) at UC Berkeley, which was a joint effort by [Daniel Geng](#), [Ryan Tabrizi](#), and [Hang Gao](#), advised by [Liyue Shen](#), [Andrew Owens](#), and [Alexei Efros](#).

Contents

- [Starter code](#)
- [Part 1](#): Training a Single-Step Denoising UNet
- [Part 2](#): Training a Diffusion Model
- [Extra credit](#)
- [Submission instructions](#)

Starter Code

Download the starter code [here](#).

Part 1: Training a Single-Step Denoising UNet

Let's warm up by building a simple one-step denoiser. Given a noisy image z , we aim to train a denoiser D_θ such that it maps z to a clean image x . To do so, we can optimize over an L2 loss:

$$L = \mathbb{E}_{z,x} \|D_\theta(z) - x\|^2 \quad (\text{B})$$

1.1 Implementing the UNet

In this project, we implement the denoiser as a [UNet](#). It consists of a few downsampling and upsampling blocks with skip connections.

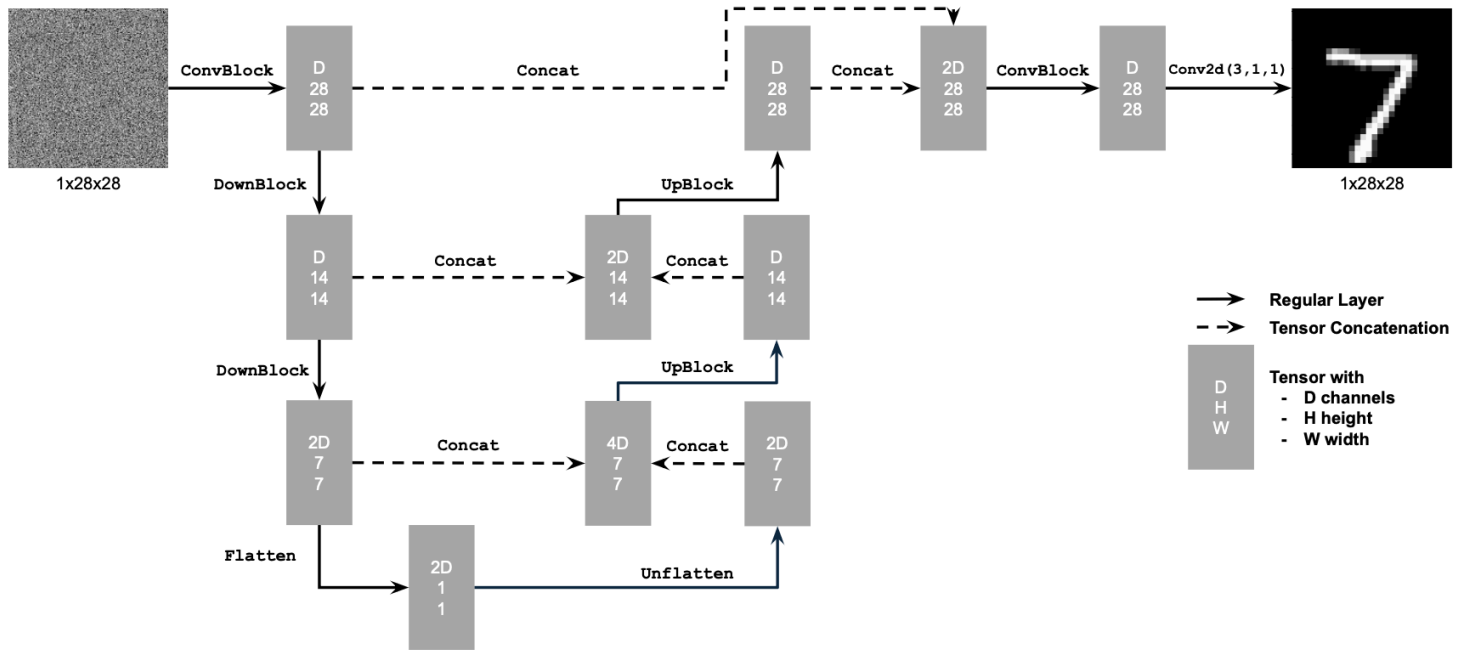


Figure 1: Unconditional UNet

The diagram above uses a number of standard tensor operations defined as follows:

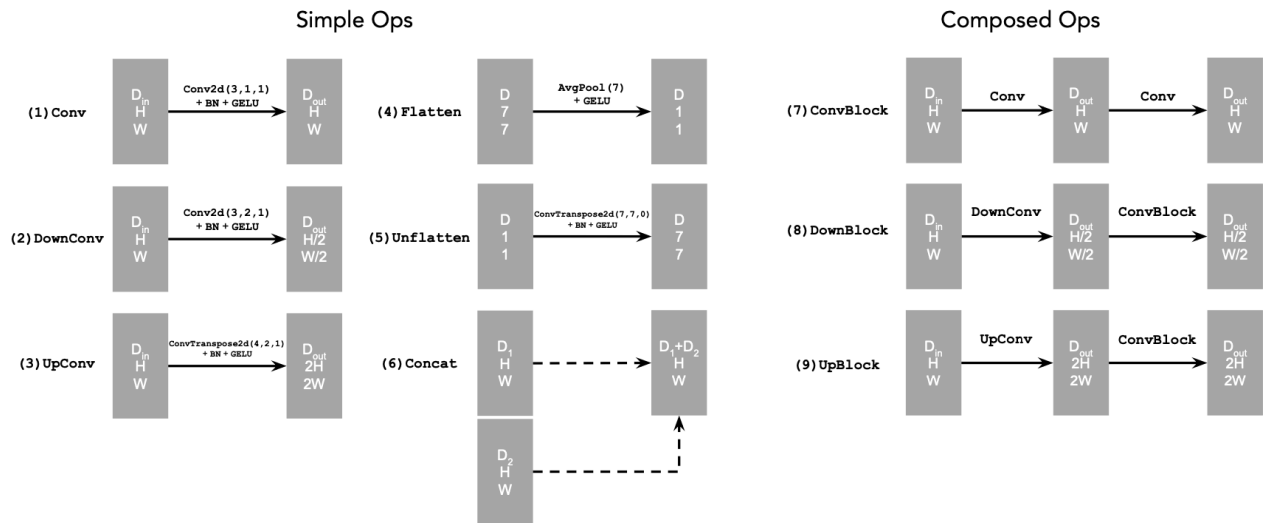


Figure 2: Standard UNet Operations

where:

- **Conv2d(kernel_size, stride, padding)** is `nn.Conv2d()`
- **BN** is `nn.BatchNorm2d()`
- **GELU** is `nn.GELU()`
- **ConvTranspose2d(kernel_size, stride, padding)** is `nn.ConvTranspose2d()`
- **AvgPool(kernel_size)** is `nn.AvgPool2d()`
- **D** is the number of hidden channels and is a hyperparameter that we will set ourselves.

At a high level, the blocks do the following:

- **(1) Conv** is a convolutional layer that doesn't change the image resolution, only the channel dimension.
- **(2) DownConv** is a convolutional layer that downsamples the tensor by 2.
- **(3) UpConv** is a convolutional layer that upsamples the tensor by 2.
- **(4) Flatten** is an average pooling layer that flattens a 7x7 tensor into a 1x1 tensor. 7 is the resulting height and width after the downsampling operations.
- **(5) Unflatten** is a convolutional layer that unflattens/upsamples a 1x1 tensor into a 7x7 tensor.
- **(6) Concat** is a channel-wise concatenation between tensors with the same 2D shape. This is simply `torch.cat()`.

We define composed operations using our simple operations in order to make our network deeper. This doesn't change the tensor's height, width, or number of channels, but simply adds more learnable parameters.

- **(7) ConvBlock**, is similar to **Conv** but includes an additional **Conv**. Note that it has the same input and output shape as **(1) Conv**.
- **(8) DownBlock**, is similar to **DownConv** but includes an additional **ConvBlock**. Note that it has the same input and output shape as **(2) DownConv**.
- **(9) UpBlock**, is similar to **UpConv** but includes an additional **ConvBlock**. Note that it has the same input and output shape as **(3) UpConv**.

1.2 Using the UNet to Train a Denoiser

Recall from equation 1 that we aim to solve the following denoising problem: Given a noisy image z , we aim to train a denoiser D_θ such that it maps z to a clean image x . To do so, we can optimize over an L2 loss

$$L = \mathbb{E}_{z,x} \|D_\theta(z) - x\|^2.$$

To train our denoiser, we need to generate training data pairs of (z, x) , where each x is a clean MNIST digit. For each training batch, we can generate z from x using the following noising process:

$$z = x + \sigma\epsilon, \quad \text{where } \epsilon \sim N(0, I). \quad (\text{B})$$

Visualize the different noising processes over $\sigma = [0.0, 0.2, 0.4, 0.5, 0.6, 0.8, 1.0]$, assuming normalized $x \in [0, 1]$. It should be similar to the following plot:

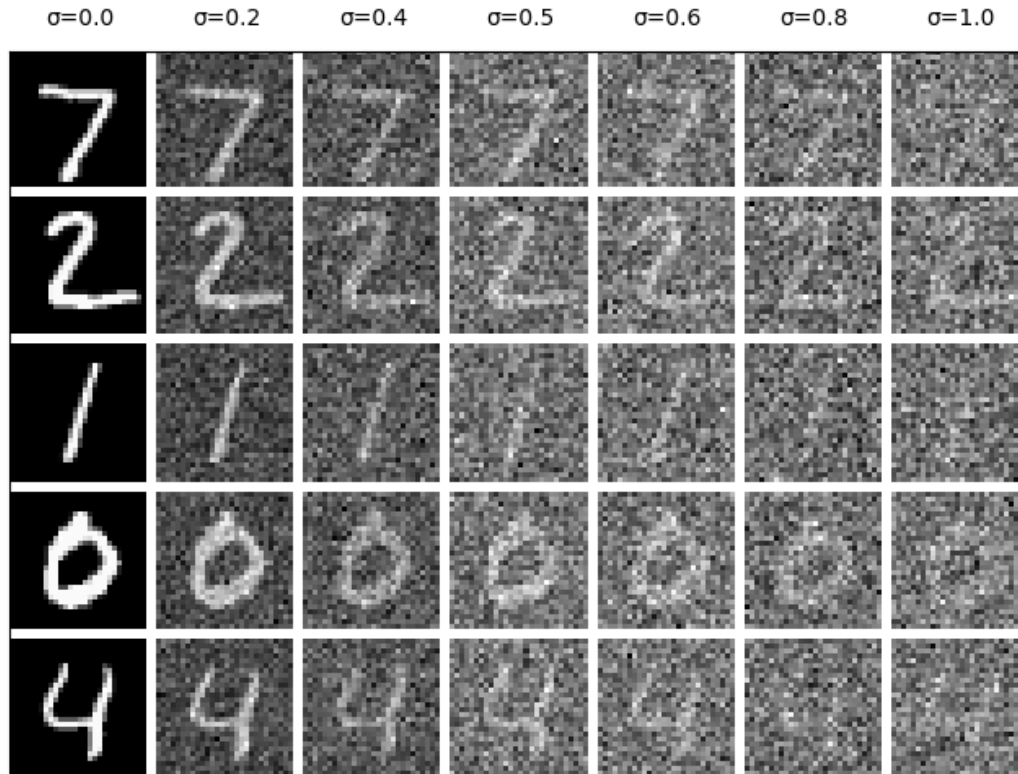


Figure 3: Varying levels of noise on MNIST digits

1.2.1 Training

Now, we will train the model to perform denoising.

- **Objective:** Train a denoiser to denoise noisy image z with $\sigma = 0.5$ applied to a clean image x .
- **Dataset and dataloader:** Use the MNIST dataset via `torchvision.datasets.MNIST` with flags to access training and test sets. Train only on the training set. Shuffle the dataset before creating the dataloader. Recommended batch size: 256. We'll train over our dataset for 5 epochs.
 - You should only noise the image batches when fetched from the dataloader so that in every epoch the network will see new noised images, improving generalization.
- **Model:** Use the UNet architecture defined in section 1.1 with recommended hidden dimension $D = 128$.
- **Optimizer:** Use Adam optimizer with learning rate of $1e-4$.

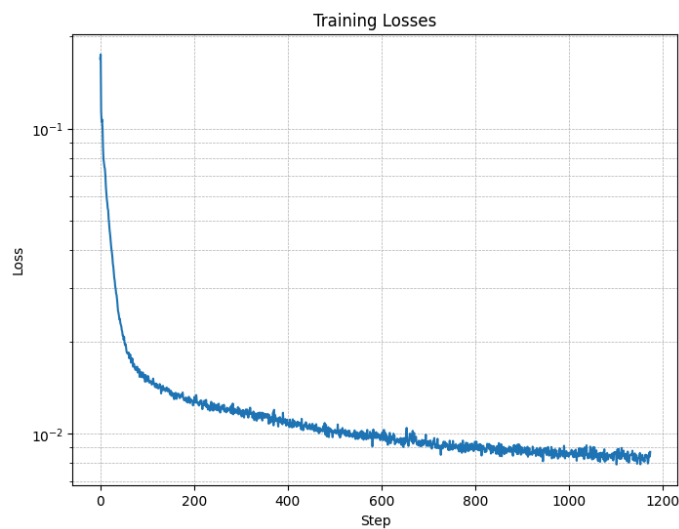


Figure 4: Training Loss Curve

You should visualize denoised results on the test set at the end of training. Display sample results after the 1st and 5th epoch.

They should look something like these:

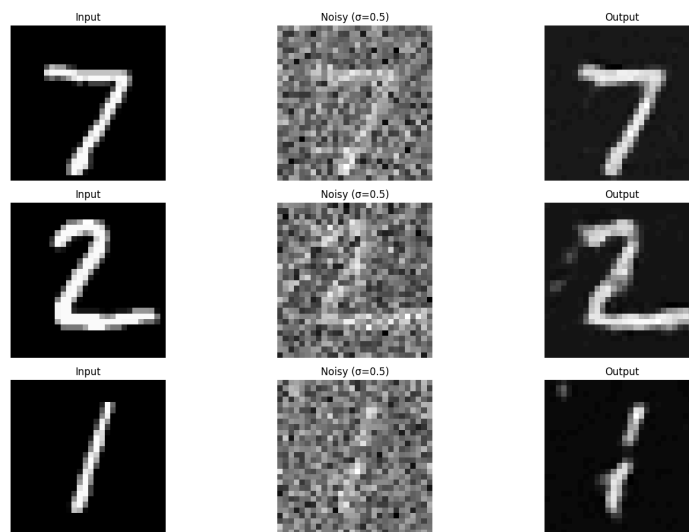


Figure 5: Results on digits from the test set after 1 epoch of training

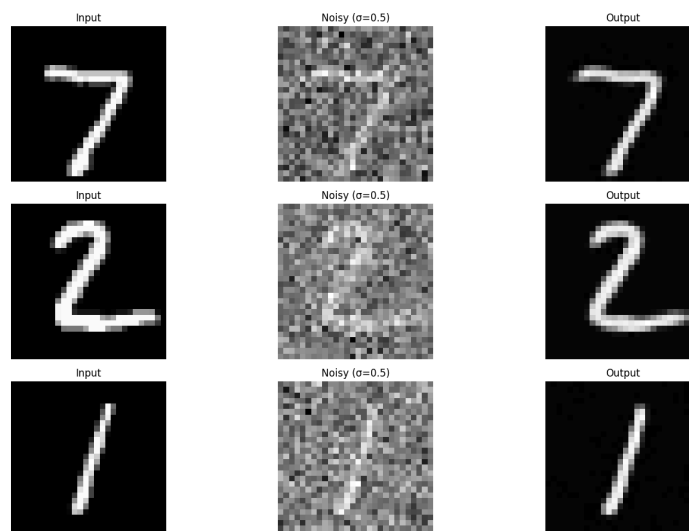


Figure 6: Results on digits from the test set after 5 epochs of training

1.2.2 Out-of-Distribution Testing

Our denoiser was trained on MNIST digits noised with $\sigma = 0.5$. Let's see how the denoiser performs on different σ 's that it wasn't trained for.

Visualize the denoiser results on test set digits with varying levels of noise $\sigma = [0.0, 0.2, 0.4, 0.5, 0.6, 0.8, 1.0]$.

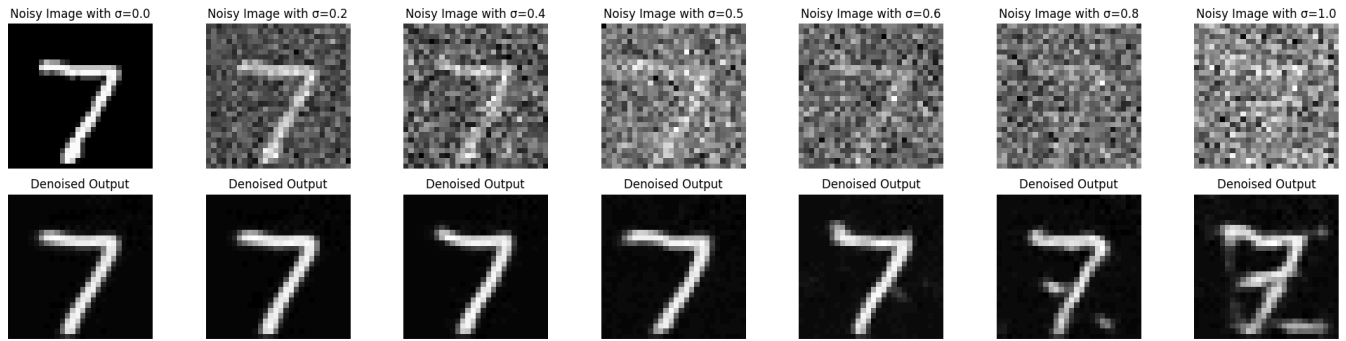


Figure 7: Results on digits from the test set with varying noise levels.

Part 2: Training a Diffusion Model

Now, we are ready for diffusion, where we will train a UNet model that can iteratively denoise an image. We will implement [DDPM](#) in this part.

Let's revisit the problem we solved in equation B.1:

$$L = \mathbb{E}_{z,x} \|D_{\theta}(z) - x\|^2.$$

We will first introduce one small difference: we can change our UNet to predict the added noise ϵ instead of the clean image x . Mathematically, these are equivalent since $x = z - \sigma\epsilon$ (equation B.2). Therefore, we can turn equation B.1 into the following:

$$L = \mathbb{E}_{\epsilon,z} \|\epsilon_{\theta}(z) - \epsilon\|^2 \quad (\text{B})$$

where ϵ_{θ} is a UNet trained to predict noise.

For diffusion, we eventually want to sample a pure noise image $\epsilon \sim N(0, I)$ and generate a realistic image x from the noise. However, one-step denoising does not yield good results. Instead, we need to *iteratively* denoise the image for better results.

We can use the following equation to generate noisy images x_t from x_0 for some timestep t for $t \in \{0, 1, \dots, T\}$:

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \quad \text{where } \epsilon \sim N(0, 1).$$

Intuitively, when $t = 0$ we want x_t to be the clean image x_0 , when $t = T$ we want x_t to be pure noise ϵ , and for $t \in \{1, \dots, T-1\}$, x_t should be some linear combination of the two. The precise derivation of $\bar{\alpha}$ is beyond the scope of this project (see [DDPM paper](#) for more details). Here, we provide you with the DDPM recipe to build a list $\bar{\alpha}$ for $t \in \{0, 1, \dots, T\}$ utilizing lists α and β :

- Create a list β of length T such that $\beta_0 = 0.0001$ and $\beta_T = 0.02$ and all other elements β_t for $t \in \{1, \dots, T-1\}$ are evenly spaced between the two.
- $\alpha_t = 1 - \beta_t$
- $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ is a cumulative product of α_s for $s \in \{1, \dots, t\}$.

Because we are working with simple MNIST digits, we can set T to 300, a relatively small number. Observe how $\bar{\alpha}_t$ is close to 1 for small t and close to 0 for T . β is known the variance schedule; it controls the amount of noise added at each timestep.

Now, to denoise image x_t , we could simply apply our UNet ϵ_{θ} on x_t and get the noise ϵ . However, this won't work very well because the UNet is expecting the noisy image have a noise variance $\sigma = 0.5$ for best results, but the variance of x_t varies with t . One could train T separate UNets, but it is much easier to simply condition a single UNet with timestep t , giving us our final objective:

$$L = \mathbb{E}_{\epsilon, x_0, t} \|\epsilon_{\theta}(x_t, t) - \epsilon\|^2. \quad (\text{B})$$

2.1 Adding Time Conditioning to UNet

We need a way to inject scalar t into our UNet model to condition it. There are many ways to do this. Here is what we suggest:

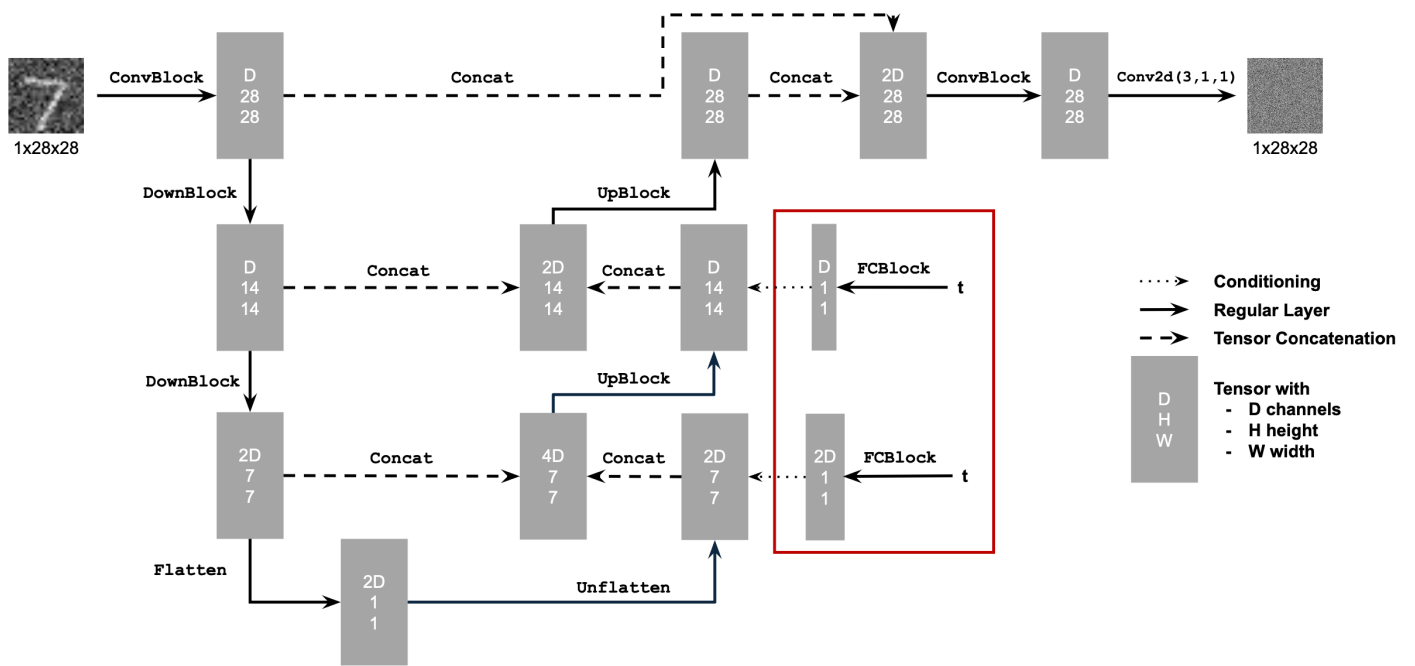


Figure 8: Conditioned UNet

This uses a new operator called **FCBlock** (fully-connected block) which we use to inject the conditioning signal into the UNet:

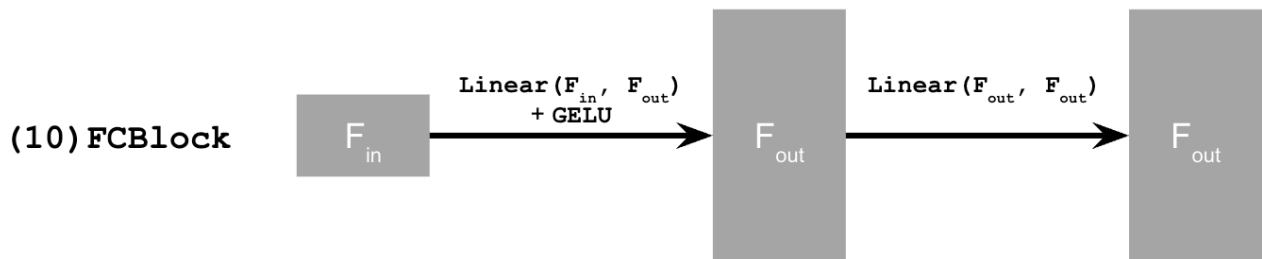


Figure 9: FCBlock for conditioning

Here **Linear(F_in, F_out)** is a linear layer with **F_in** input features and **F_out** output features. You can implement it using `nn.Linear`.

Since our conditioning signal t is a scalar, **F_in** should be of size 1. We also recommend that you normalize t to be in the range $[0, 1]$ before embedding it, i.e. pass in $\frac{t}{T}$.

You can embed t by following this pseudo code:

```
fc1_t = FCBlock(...)
fc2_t = FCBlock(...)

# the t passed in here should be normalized to be in the range [0, 1]
t1 = fc1_t(t)
t2 = fc2_t(t)

# Follow diagram to get unflatten.
# Replace the original unflatten with modulated unflatten.
unflatten = unflatten + t1
# Follow diagram to get up1.
...
# Replace the original up1 with modulated up1.
up1 = up1 + t2
# Follow diagram to get the output.
...
```

2.2 Training the UNet

Training our time-conditioned UNet $\epsilon_\theta(x_t, t)$ is now pretty easy. Basically, we pick a random image from the training set, a random t , and train the denoiser to predict the noise in x_t . We repeat this for different images and different t values until the model converges and we are happy.

```
Algorithm 1 Training
1: Precompute  $\bar{\sigma}$ 
2: repeat
3:    $x_0 \sim$  clean image from training set
4:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
5:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
6:    $x_t = \sqrt{\bar{\sigma}_t} x_0 + \sqrt{1 - \bar{\sigma}_t} \epsilon$ 
7:    $\hat{\epsilon} = \epsilon_\theta(x_t, t)$ 
8:   Take gradient descent step on  $\|\hat{\epsilon} - \epsilon\|^2$ 
9: until happy
```

Algorithm B.1. Training time-conditioned UNet

- **Objective:** Train a time-conditioned UNet $\epsilon_\theta(x_t, t)$ to predict the noise in x_t given a noisy image x_t and a timestep t .
- **Dataset and dataloader:** Use the MNIST dataset via `torchvision.datasets.MNIST` with flags to access training and test sets. Train only on the training set. Shuffle the dataset before creating the dataloader. Recommended batch size: 128. We'll train over our dataset for 20 epochs.
 - As shown in algorithm B.1, You should only noise the image batches when fetched from the dataloader.
- **Model:** Use the time-conditioned UNet architecture defined in section 2.1 with recommended hidden dimension $D = 64$. Follow the diagram and pseudocode for how to inject the conditioning signal t into the UNet. Remember to normalize t before embedding it.
- **Optimizer:** Use Adam optimizer with an initial learning rate of $1e-3$. We will be using an exponential learning rate decay scheduler with a gamma of $0.1^{(1.0/\text{num_epoch})}$. This can be implemented using `scheduler = torch.optim.lr_scheduler.ExponentialLR(...)`. You should call `scheduler.step()` after every epoch.

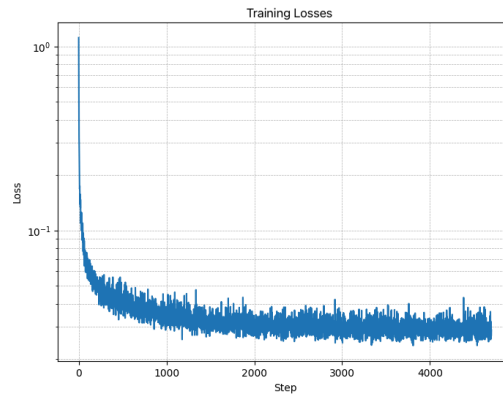


Figure 10: Time-Conditioned UNet training loss curve

2.3 Sampling from the UNet

Now that we trained the UNet, we can do sampling using the process below.

Algorithm 2 Sampling

```

1: Precompute  $\beta_t, \alpha_t$  and  $\bar{\alpha}_t$ 
2:  $x_T \sim \mathcal{N}(0, I)$ 
3: for  $t$  from  $T$  to 1, step size  $-1$  do
4:    $z \sim \mathcal{N}(0, I)$  if  $t > 1$ , else  $z = 0$ 
5:    $x_t = \frac{\alpha_t}{\sqrt{\beta_t}}(x_{t+1} - \sqrt{1 - \alpha_t/\beta_t}\epsilon_\theta(x_{t+1}, t)) + \sqrt{\beta_t}z$ 
6:    $x_{t-1} = \frac{\sqrt{\alpha_{t-1}}}{1 - \alpha_t}x_t + \frac{\sqrt{\alpha_t(1 - \alpha_{t-1})}}{1 - \alpha_t}x_t + \sqrt{\beta_t}z$ 
7: end for
8: return  $x_0$ 

```

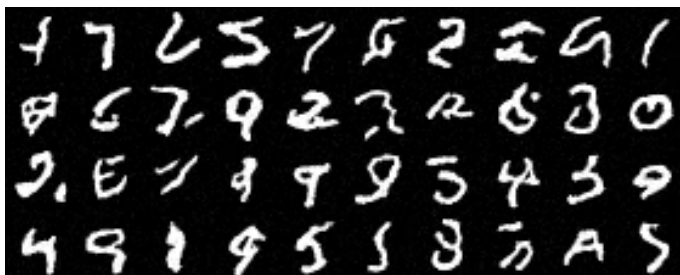
Algorithm B.2. Sampling from time-conditioned UNet



Epoch 1



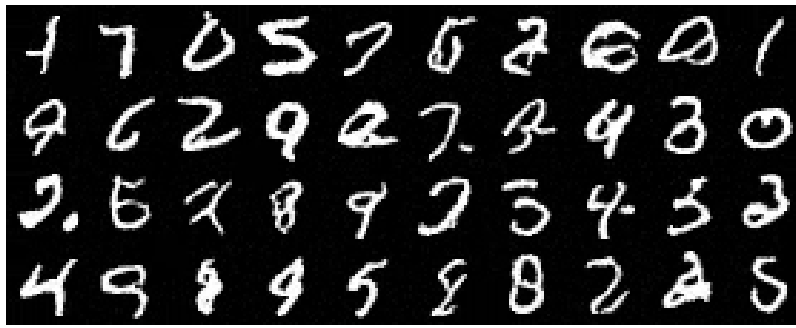
Epoch 5



Epoch 10



Epoch 15



Epoch 20

2.4 Adding Class-Conditioning to UNet

To make the results better and give us more control for image generation, we can also optionally condition our UNet on the class of the digit 0-9. This will require adding 2 more **FCBlocks** to our UNet but, we suggest that for class-conditioning vector c , you make it a one-hot vector instead of a single scalar. Because we still want our UNet to work without it being conditioned on the class, we implement dropout where 10% of the time ($p_{\text{uncond}} = 0.1$) we drop the class conditioning vector c by setting it to 0. Here one way to condition our UNet $\epsilon_{\theta}(x_t, t, c)$ on both time t and class c :

```
fc1 t = FCBlock(...)
fc1 c = FCBlock(...)
fc2 t = FCBlock(...)
fc2_c = FCBlock(...)

t1 = fc1 t(t)
c1 = fc1 c(c)
t2 = fc2 t(t)
c2 = fc2_c(c)

# Follow diagram to get unflatten.
# Replace the original unflatten with modulated unflatten.
unflatten = c1 * unflatten + t1
# Follow diagram to get up1.
...
# Replace the original up1 with modulated up1.
up1 = c2 * up1 + t2
# Follow diagram to get the output.
...
```

Training for this section will be the same as time-only, with the only difference being the conditioning vector c and doing unconditional generation periodically.

Algorithm 3 Class-Conditioned Training

- 1: Precompute $\bar{\alpha}$
- 2: repeat
- 3: $x_0, c \sim$ clean image and label from training set
- 4: Make c into a one-hot vector
- 5: (with probability p_{uncond} set c to zero-vector).
- 6: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 7: $\epsilon \sim \mathcal{N}(0, I)$
- 8: $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$
- 9: $\hat{\epsilon} = \epsilon_{\theta}(x_t, t, c)$
- 10: Take gradient descent step on $\nabla_{\theta} \| \epsilon - \hat{\epsilon} \|^2$
- 11: until happy

Algorithm B.3. Training class-conditioned UNet

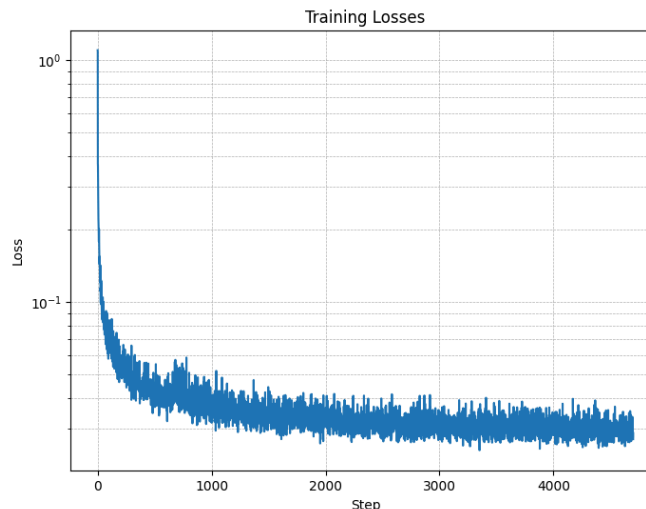


Figure 11: Class-conditioned UNet training loss curve

2.5 Sampling from the Class-Conditioned UNet

Follow the sampling process below and use classifier-free guidance with $\gamma = 5.0$.

Algorithm 4 Class-Conditioned Sampling

```

1: input: one-hot vector  $c$ , classifier-free guidance scale  $\gamma$ 
2:  $x_T \sim \mathcal{N}(0, I)$ 
3: for  $t$  from  $T$  to 1, step size  $-1$  do
4:    $z \sim \mathcal{N}(0, I)$ ; if  $t > 1$ , else  $z = 0$ 
5:    $e_u = e_\theta(x_t, t, 0)$ 
6:    $e_c = e_\theta(x_t, t, c)$ 
7:    $\epsilon = e_u + \gamma(e_c - e_u)$  ▷ Classifier-free guidance
8:    $\hat{x}_0 = -\frac{\epsilon}{\alpha_t} (x_t - \sqrt{1 - \alpha_t} z)$ 
9:    $x_{t-1} = \frac{\sqrt{\alpha_{t-1}}}{1 - \alpha_t} \hat{x}_0 + \frac{\sqrt{\alpha_t(1 - \alpha_{t-1})}}{1 - \alpha_t} x_t + \sqrt{\beta_t} \#$ 
10: end for
11: return  $x_0$ 

```

Algorithm B.4. Sampling from class-conditioned UNet



Epoch 1



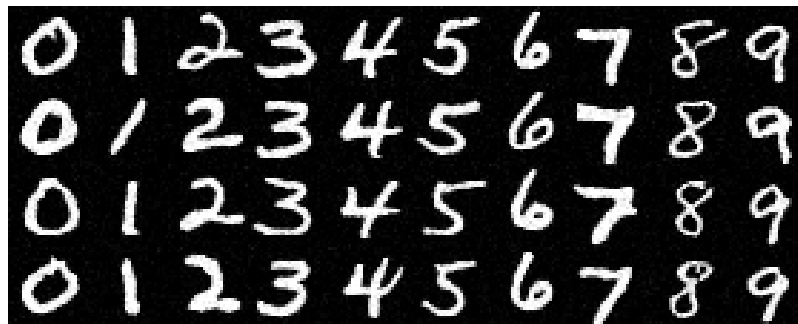
Epoch 5



Epoch 10



Epoch 15



Epoch 20

Extra Credit

- Create your own sampling gifs similar to the ones shown above.
- Train a digit classifier and implement classifier-guided sampling. Compare your results to classifier-free guidance.
- Attempt to train a DDPM on a bigger and more complicated dataset, such as CIFAR-100 or [Imagenette](#).
- Improve the network architecture for time-conditional generation. For ease of explanation and implementation, our network architecture above is pretty simple. Modify the network (e.g., with skip connections or self-attention) and the time embedding so you can sample even better results. Compare the results between the original and modified networks in your report, and discuss the changes you made and any observed differences.
- Implement a more advanced model such as Rectified Flow.
 - Implement [rectified flow](#), which is the state of art diffusion model.
 - You can reference any code on github, but your implementation needs to follow the same code structure as our DDPM implementation.
 - In other words, the code change required should be minimal: only changing the forward and sample functions.

Submission Instructions:

The assignment deliverables are as follows. If you are working in a pair, only one designated student should make the submission.

1. Upload the following files to [Canvas](#):
 1. **netid_mp5_output.pdf**: Your IPython notebook with output cells converted to **PDF format**

2. **mp5.ipynb**: Your top-level notebook
3. **netid_mp5_report.pdf**: Your assignment report (using [this template](#)) in PDF format
4. **(Optional)** Any additional files related to extra credit

Please refer to [course policies](#) on collaborations, late submission, etc.