

IFT6135A – W2025

## Assignment 3: Practical Part

Université de Montréal

Felix Wilhelmy

October 4, 2025

---

**Due Date: May 2, 23:00**

Instructions

- For all questions that are not graded only on the answer, show your work! Any problem without work shown will get no marks regardless of the correctness of the final answer.
- Please try to use a document preparation system such as LaTeX. *If you write your answers by hand, note that you risk losing marks if your writing is illegible without any possibility of regrade, at the discretion of the grader.*
- Submit your answers electronically via the course GradeScope. *Incorrectly assigned answers can be given 0 automatically at the discretion of the grader.* To assign answers properly, please:
  - Make sure that the top of the first assigned page is the question being graded.
  - Do not include any part of answer to any other questions within the assigned pages.
  - Assigned pages need to be placed in order.
  - For questions with multiple parts, the answers should be written in order of the parts within the question.
- In the code, each part to fill is referenced by a TODO and ‘Not Implemented Error’
- Questions requiring written responses should be short and concise when necessary. Unnecessary wordiness will be penalized at the grader’s discretion.
- Please sign the agreement below.
- It is your responsibility to follow updates to the assignment after release. All changes will be visible on Overleaf and Piazza.
- Any questions should be directed towards the TAs for this assignment: *Vitória Barin Pacela, Philippe Martin.*

For this assignment, the GitHub link is the following: [https://github.com/philmar1/Teaching\\_IFT6135---Assignment-3---H25](https://github.com/philmar1/Teaching_IFT6135---Assignment-3---H25)

**I acknowledge I have read the above instructions and will abide by them throughout this assignment. I further acknowledge that any assignment submitted without the following form completed will result in no marks being given for this portion of the assignment.**

Name: **Félix Wilhelmy**

UdeM Student ID: **20333575**

Signature: \_\_\_\_\_

---

## Introduction

All experiments were conducted on Google Colab utilizing T4 GPUs. All models were trained for 20 epochs; however, due to a saving bug, checkpoints for DDPM and CFG models at epoch 20 were not properly saved. Consequently, the results presented for these models correspond to epoch 19.

---

## Question 1 : Variational Autoencoder (VAE)

### Question 1.6 : Training & Validation Loss

The loss curves for both training and validation exhibit a steady decline over the 20-epoch run, ultimately converging to final values of 103.8882 (training) and 103.6296 (validation), thereby meeting the prescribed target of 104. As shown in Figure 1, the two curves track closely throughout, indicating stable optimization without significant overfitting. The plateau observed toward the final epochs further confirms consistent learning dynamics across both datasets, with no divergence suggesting under- or overfitting.

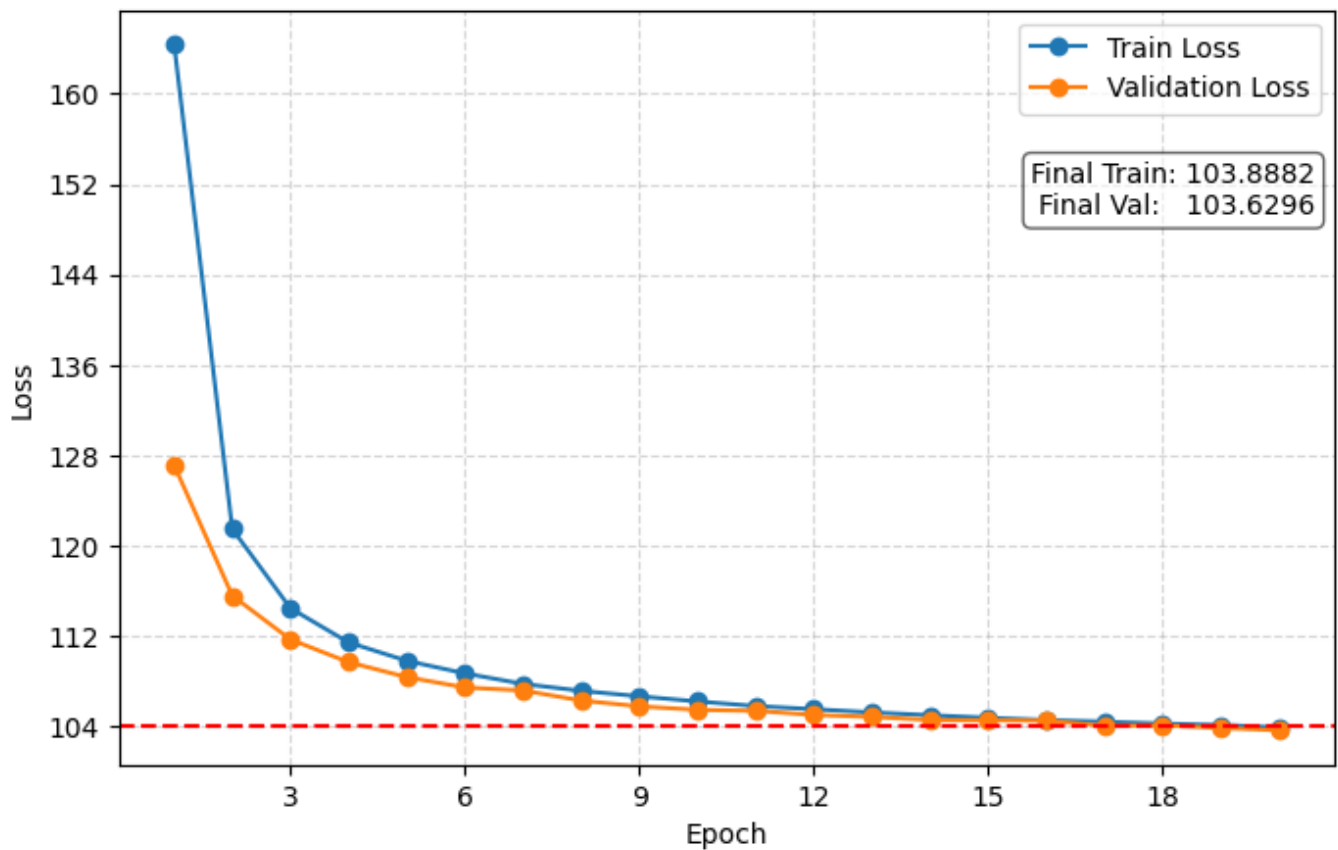


Figure 1: Training and validation loss curves over 20 epochs for the MNIST VAE. The dashed red line marks the target validation loss of 104, and the final loss values are annotated in the top-right corner.

---

## Code for the generation of the figure

```
1 def plot_loss(history_path: str, save_path: str = None, show: bool = False):
2     if not os.path.isfile(history_path):
3         raise FileNotFoundError(f"No file found at {history_path!r}")
4
5     with open(history_path, 'r') as fp:
6         history = json.load(fp)
7
8     train_loss = history.get('train_loss', [])
9     val_loss = history.get('val_loss', [])
10
11     num_epochs = len(train_loss)
12     epochs = range(1, num_epochs + 1)
13
14     fig, ax = plt.subplots(figsize=(8, 5))
15     ax.plot(epochs, train_loss, marker='o', label='Train Loss')
16     ax.plot(epochs, val_loss, marker='o', label='Validation Loss')
17
18     # Horizontal red line at y=104
19     ax.axhline(y=104, color='red', linestyle='--')
20
21     ax.yaxis.set_major_locator(MaxNLocator(integer=True))
22     ax.set_xticks(range(1, num_epochs + 1)) # Ticks for each epoch
23     ax.xaxis.set_major_locator(MaxNLocator(integer=True))
24
25     # Annote final loss in the top-right corner
26     final_tr = train_loss[-1]
27     final_va = val_loss[-1]
28     text = f"Final Train: {final_tr:.4f}\nFinal Val: {final_va:.4f}"
29     ax.text(
30         0.985, 0.8, text,
31         transform=ax.transAxes,
32         fontsize=10,
33         ha='right', va='top',
34         bbox=dict(boxstyle='round,pad=0.3', facecolor='white', alpha=0.6)
35     )
36
37     ax.set_xlabel('Epoch')
38     ax.set_ylabel('Loss')
39     ax.grid(True, linestyle='--', alpha=0.5)
40     ax.legend()
41
42     if save_path:
43         plt.savefig(save_path, bbox_inches='tight')
44         print(f"Saved loss plot to {save_path}")
45     if show:
46         plt.show()
47     plt.close(fig)
```

---

### Question 1.7 : Sample Generation

The trained VAE produces samples that, while somewhat blurred, remain identifiable as MNIST digits and cover a variety of classes, showing no evidence of mode collapse. As illustrated in Figure 2, the digits exhibit thicker or more faded strokes and lack the sharpness of DDPM and CFG outputs. Such blurriness arises from the VAE decoder’s implicit averaging under uncertainty, yet the overall coherence and placement of strokes confirm that the model has captured the MNIST manifold and learned to generate plausible handwriting patterns.



Figure 2: An  $8 \times 8$  grid of 64 images sampled from the trained VAE decoder (epoch 20), demonstrating the model’s generative capacity on MNIST.

Despite that, the digits are generally coherent and fall within the manifold of realistic MNIST examples, showing that the VAE has learned to generate plausible handwriting samples (albeit with lower sharpness).

---

## Code for the generation of the figure

```
1 def generate_images(  
2     model: VAE,  
3     num_images: int = 64,  
4     device: torch.device = torch.device('cpu'),  
5     show: bool = False,  
6     file_path: Optional[str] = None  
7 ) -> torch.Tensor:  
8     model.eval()  
9     model.to(device)  
10  
11     with torch.no_grad():  
12         latent_dim = 20 # Hardcoded for our model  
13         z = torch.randn(num_images, latent_dim, device=device)  
14         x_flat = model.decode(z) # (num_images, 784)  
15         x = x_flat.view(num_images, 1, 28, 28) # (num_images, 1, 28, 28)  
16  
17         if show or file_path:  
18             grid = torchvision.utils.make_grid(x, nrow=8, padding=2)  
19             plt.figure(figsize=(6, 6))  
20             plt.imshow(grid.permute(1, 2, 0).cpu(), cmap='gray')  
21             plt.axis('off')  
22             if file_path:  
23                 plt.savefig(file_path, bbox_inches='tight')  
24             if show:  
25                 plt.show()  
26             plt.close()  
27  
28     return x
```

---

## Question 1.8 : Latent Traversals

To assess disentanglement, we perform two sets of traversals over all 20 latent dimensions: one starting from a random prior sample  $z_0 \sim \mathcal{N}(0, I)$ , and another from the encoding of an observed MNIST image. For each latent factor  $i \in \{0, \dots, 19\}$ , we perturb the base vector by  $-2\varepsilon$ ,  $-\varepsilon$ ,  $0$ ,  $+\varepsilon$ , and  $+2\varepsilon$  (with  $\varepsilon = 2.0$ ) and decode the results. Each row in Figures 3–4 corresponds to one latent dimension; columns show increasing offsets.

The results do not reveal strongly disentangled factors in the VAE’s learned representation. Varying each latent dimension (while fixing others) causes changes in the generated digit images, but these changes are entangled. For example, adjusting a particular latent dimension might simultaneously alter a digit’s stroke thickness and its shape or orientation, rather than isolating one property alone. Nevertheless, a few factors show relatively isolated effects:

- **Factor 15:** Consistently shifts decoded outputs between “0”-like shapes when negatively perturbed and “1”-like shapes when positively perturbed.
- **Factor 13:** Primarily modulates stroke width across digits.
- **Factor 17:** When increased, tends to morph various digits toward the appearance of a “4.”

Overall, this indicates that even if some factor seem to have obtained some disentanglement, the latent space of the VAE is largely entangled: the generative factors (such as style, stroke width, or positioning) are mostly distributed across multiple latent dimensions rather than being neatly separated one per dimension.



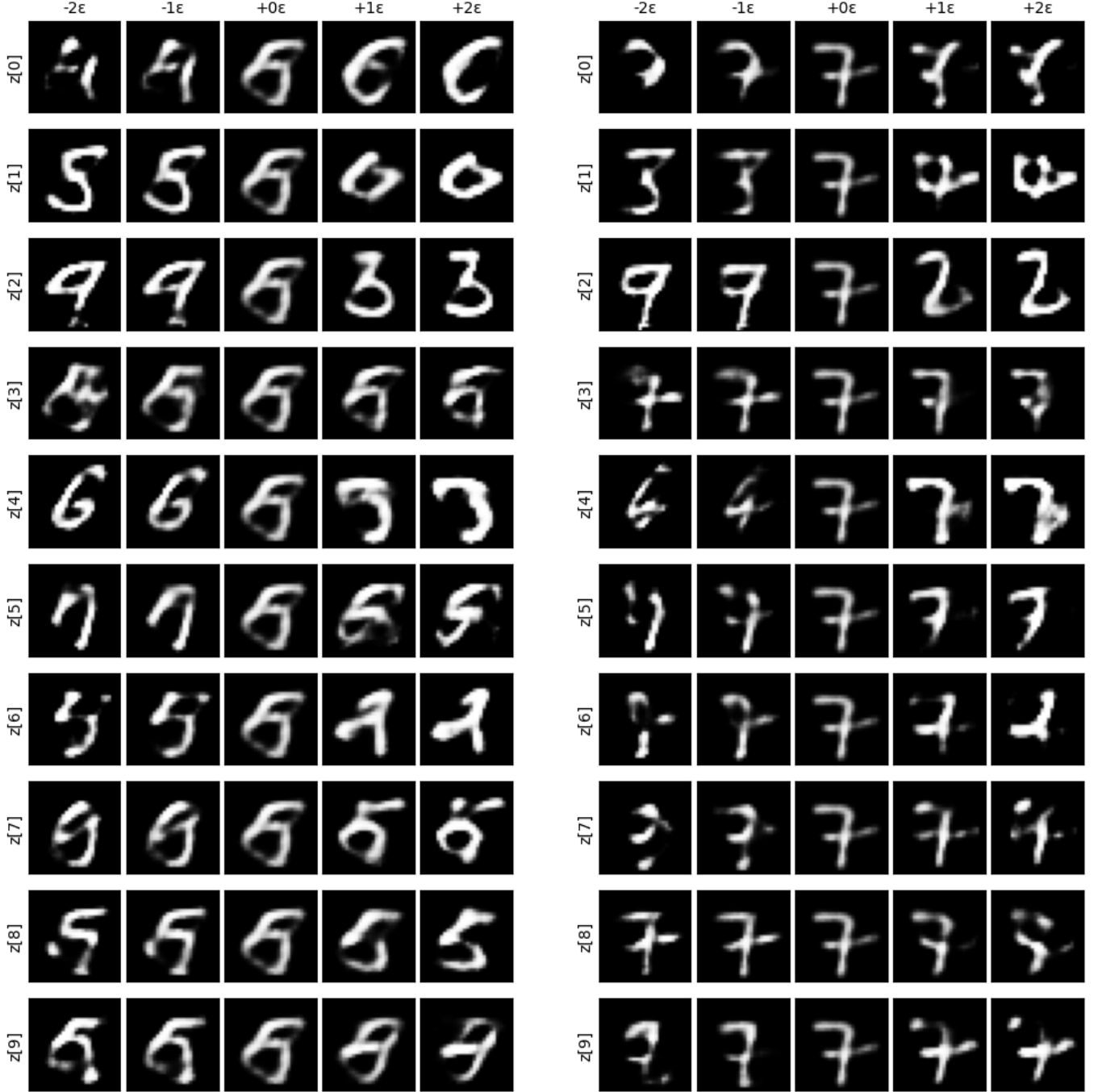


Figure 3: Latent dimensions 0–9. **Left:** traversal from a Gaussian prior sample. **Right:** traversal from an encoded MNIST image. Each row  $i$  shows perturbations of  $z_i$  by  $\{-2, -1, 0, +1, +2\}\epsilon$ .



Figure 4: Latent dimensions 10–19. **Left:** traversal from a Gaussian prior sample. **Right:** traversal from an encoded MNIST image. Each row  $i$  shows perturbations of  $z_i$  by  $\{-2, -1, 0, +1, +2\}\epsilon$ .

---

## Code for the latent traversals

```
1 def generate_latent_traversals(  
2     model: VAE,  
3     n_steps: int = 5,  
4     eps: float = 2.0,  
5     device: torch.device = torch.device('cpu'),  
6     data_sample: Optional[torch.Tensor] = None,  
7 ) -> torch.Tensor:  
8     model.eval()  
9     model.to(device)  
10  
11     latent_dim = 20 # Hardcoded for our model  
12     if data_sample is not None: # Sampled from X (MNIST)  
13         x_flat = data_sample.to(device).view(1, -1)  
14         mu, logvar = model.encode(x_flat)  
15         z0 = model.reparameterize(mu, logvar)  
16     else: # Sampled from the prior  $Z = N(0, I)$   
17         z0 = torch.randn(1, latent_dim, device=device)  
18  
19     traversals = []  
20     half_span = (n_steps - 1) / 2  
21  
22     with torch.no_grad():  
23         for dim in range(latent_dim):  
24             row = []  
25             for step in range(n_steps):  
26                 z = z0.clone()  
27                 # Perturb the dim-th latent variable  
28                 z[0, dim] += (step - half_span) * eps  
29                 x = model.decode(z) # (1, P), P = H*W  
30                 side = int(math.sqrt(x.size(1)))  
31                 row.append(x.view(1, 1, side, side)) # (1, 1, H, W)  
32                 # Stack steps → (n_steps, 1, H, W)  
33                 traversals.append(torch.cat(row, dim=0))  
34  
35     # Stack dims → (n_factors, n_steps, 1, H, W)  
36     return torch.stack(traversals, dim=0).cpu()
```

## Code for the generation of the figure

```
1 def plot_traversals(  
2     traversals: torch.Tensor,  
3     save_path: Optional[str] = None,  
4     show: bool = True,  
5 ):  
6     n_factors, n_steps, _, H, W = traversals.shape
```

---

```

7  center = n_steps // 2
8  mid_x = (W - 1) / 2
9  mid_y = (H - 1) / 2
10
11  # Iterate in blocks of 10 latent dims
12  for block_start in range(0, n_factors, 10):
13      block_end = min(block_start + 10, n_factors)
14      block_size = block_end - block_start
15
16      fig, axes = plt.subplots(
17          block_size, n_steps,
18          figsize=(4.8, 10),
19          gridspec_kw={'wspace': 0},
20          squeeze=False
21      )
22
23      for row, factor in enumerate(range(block_start, block_end)):
24          for step in range(n_steps):
25              ax = axes[row, step]
26              ax.imshow(traversals[factor, step, 0], cmap='gray', vmin=0, vmax=1)
27
28              # Hide ticks and spines
29              ax.set_xticks([])
30              ax.set_yticks([])
31
32              # Top row x-tick labels
33              if row == 0:
34                  ax.xaxis.set_ticks([mid_x])
35                  ax.xaxis.set_ticklabels([f'{step - center:+d}'], fontsize=10)
36                  ax.xaxis.tick_top()
37                  ax.tick_params(axis='x', length=0)
38
39              # First collumn y-tick labels
40              if step == 0:
41                  ax.yaxis.set_ticks([mid_y])
42                  ax.yaxis.set_ticklabels([f'z[{factor}]'], fontsize=10, rotation=90, va='center')
43                  ax.tick_params(axis='y', length=0)
44
45      plt.tight_layout()
46      if save_path:
47          # Append block index to filename before extension
48          base, ext = os.path.splitext(save_path)
49          part_path = f"{base}_part{block_start//10}-{ext}"
50          os.makedirs(os.path.dirname(part_path), exist_ok=True)
51          fig.savefig(part_path, bbox_inches='tight')
52
53      if show:
54          plt.show()
55
56      plt.close(fig)

```

---

### Question 1.9 : Interpolation Comparison

We compare two interpolation schemes between points  $z_0, z_1$  in the VAE’s latent space and their corresponding data-space images  $x_0, x_1$ . Figure 5 (top row) shows decoded images for  $\alpha \in \{0, 0.1, \dots, 1\}$  via latent-space interpolation  $z'_\alpha = \alpha z_0 + (1 - \alpha)z_1$ . Each intermediate decode is a plausible MNIST digit that transitions smoothly from one class to the next. In contrast, Figure 5 (bottom row) displays pixel-space interpolation  $x'_\alpha = \alpha x_0 + (1 - \alpha)x_1$ , which yields ghosted or blurry superpositions that only resemble valid digits at the endpoints.

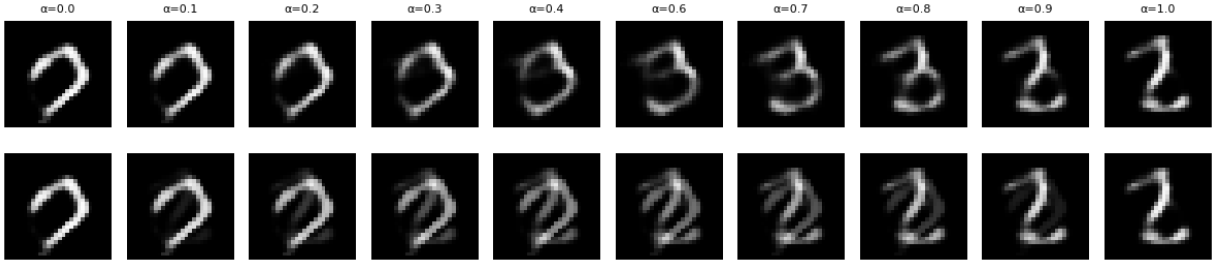


Figure 5: Comparison of interpolation methods between two random latent samples. **Top row:** latent-space interpolation yields coherent digit morphs. **Bottom row:** pixel-space interpolation produces incoherent blends except at  $\alpha = 0, 1$ .

Thus, the VAE’s latent-space interpolation produces far more coherent and meaningful intermediate outputs, highlighting that the latent space has learned a high-level representation of digit features, whereas pixel-space mixing does not respect the data’s true structure.

To further illustrate, we perform latent-space interpolations between real MNIST images. Figures 6–8 show smooth transitions learned in latent space, with digit-specific structural changes that respect the data manifold.

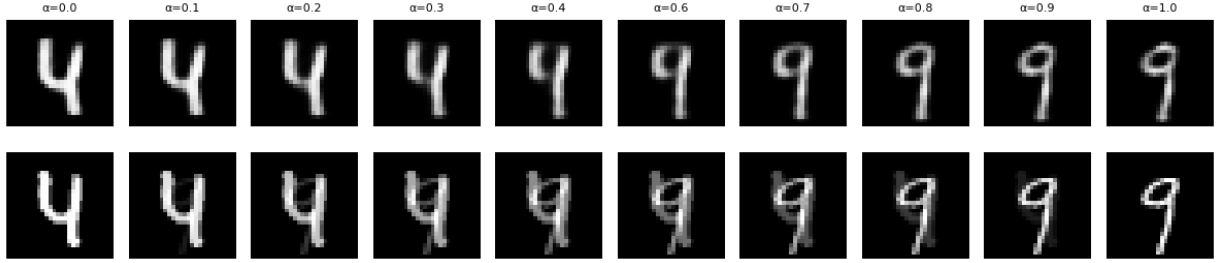


Figure 6: Latent-space interpolation between a real “4” and a real “9.” The upper loop of the “4” gradually closes to form the “9.”

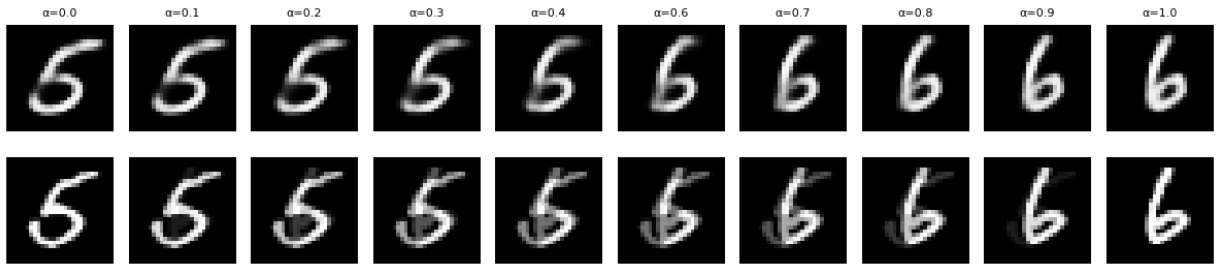


Figure 7: Latent-space interpolation between a real “5” and a real “6.” The “5” loop smoothly closes into the “6.”

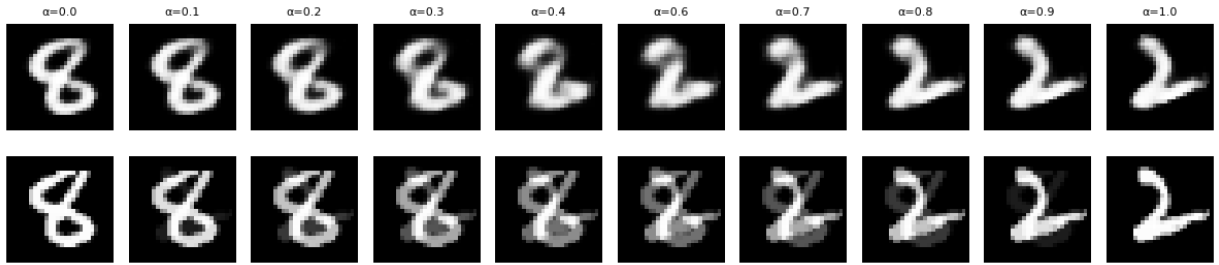


Figure 8: Latent-space interpolation between a real “8” and a real “2.” The two loops of the “8” collapse sequentially to form the “2,” though with slightly rougher transitions.

---

## Code for the generation of the figure

```
1 def plot_interpolation(
2     model: VAE,
3     endpoints: Optional[Tuple[torch.Tensor, torch.Tensor]] = None,
4     num_steps: int = 11,
5     device: torch.device = torch.device('cpu'),
6     save_path: Optional[str] = None,
7     show: bool = True
8 ):
9     model.eval()
10    model.to(device)
11
12    latent_dim = 20 # Hardcoded for our model
13
14    if endpoints is None: # Sample from prior  $Z = N(0, I)$ 
15        z0 = torch.randn(1, latent_dim, device=device)
16        z1 = torch.randn(1, latent_dim, device=device)
17        with torch.no_grad():
18            x0 = model.decode(z0).view(1,1,28,28).cpu()
19            x1 = model.decode(z1).view(1,1,28,28).cpu()
20    else: # Sample from data  $X$  (MNIST)
21        x0, x1 = endpoints
22        x0 = x0.to(device); x1 = x1.to(device)
23        flat0, flat1 = x0.view(1,-1), x1.view(1,-1)
24        with torch.no_grad(): # Encode to latent space
25            mu0, logvar0 = model.encode(flat0)
26            mu1, logvar1 = model.encode(flat1)
27            z0 = model.reparameterize(mu0, logvar0)
28            z1 = model.reparameterize(mu1, logvar1)
29    # Keep CPU copies for pixel interpolation
30    x0, x1 = x0.cpu(), x1.cpu()
31
32    # Interpolation weights
33    alphas = torch.linspace(0, 1, steps=num_steps, device=device)
34
35    # Latent-Space interpolation
36    latent_interpolation = []
37    for a in alphas:
38        z = a*z0 + (1-a)*z1
39        with torch.no_grad():
40            img = model.decode(z).view(1,1,28,28).cpu()
41        latent_interpolation.append(img)
42    latent_interpolation = torch.cat(latent_interpolation, dim=0)
43
44    # Pixel-Space interpolation
45    pixel_interpolation = []
46    for a in alphas:
47        img = a*x0 + (1-a)*x1
```

---

```

48 pixel_interpolation.append(img)
49 pixel_interpolation = torch.cat(pixel_interpolation, dim=0)
50
51 fig, axes = plt.subplots(2, num_steps, figsize=(12, 3), squeeze=False)
52 for i in range(num_steps):
53     axes[0,i].imshow(latent_interpolation[i,0], cmap='gray', vmin=0, vmax=1)
54     axes[0,i].axis('off')
55     axes[1,i].imshow(pixel_interpolation[i,0], cmap='gray', vmin=0, vmax=1)
56     axes[1,i].axis('off')
57     axes[0,i].set_title(f"={alphas[i]:.1f}", fontsize=8)
58
59 plt.tight_layout()
60 if save_path:
61     plt.savefig(save_path, bbox_inches='tight')
62 if show:
63     plt.show()
64 plt.close(fig)

```



---

## Question 2 : Unconditional Diffusion Model (DDPM)

### Question 2.1 : Sample Generation (Across Epochs)

During training, the quality of the DDPM’s generated digits improves dramatically. Early checkpoints produce samples that are noisy and often unrecognizable, but by epoch 20 the model yields well-formed MNIST digits—with only minor blurring or small artifacts remaining. As shown in Figure 9, the progressive enhancement of stroke clarity and digit structure is apparent across epochs 5, 9, 15 and 19.

A key distinction from the VAE samples is the nature of the artifacts: rather than overall blur, the diffusion outputs exhibit “shaky” lines or missing strokes. This suggests that the DDPM’s latent denoising process captures sharper, more detailed features, even if fine-grained stability (straight strokes, closed loops) can still be improved.

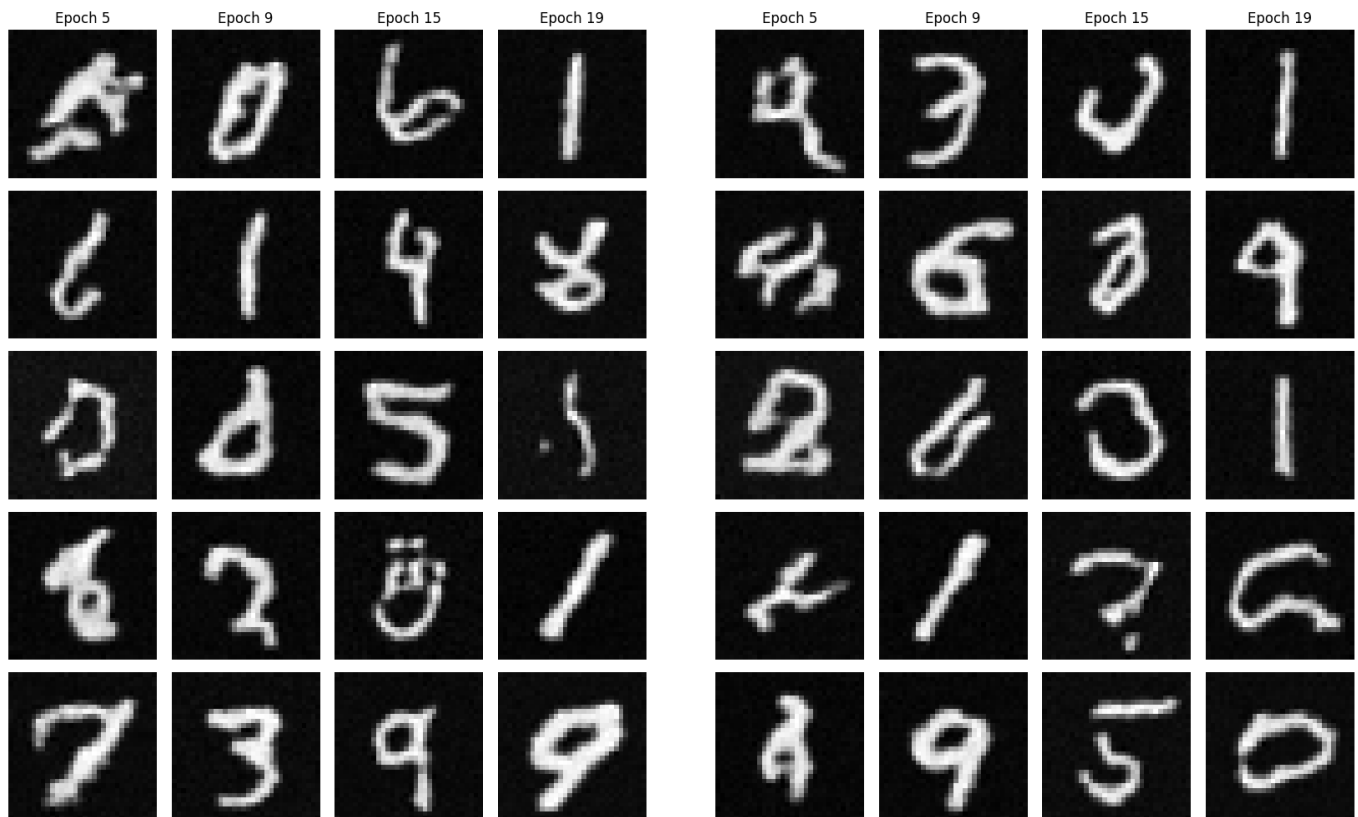


Figure 9: Unconditional DDPM samples at epochs 5, 9, 15, and 19 (columns) for four independent draws (rows). Early epochs produce noisy structures, while later epochs yield clear digit shapes.

---

Figure 10 presents a larger grid of 64 samples from the final epoch, illustrating the overall diversity and typical fidelity achieved after 20 epochs.



Figure 10: Grid of 64 MNIST images generated by the DDPM after 20 epochs, illustrating the model’s typical sample quality and diversity post-training.

To further enhance sample quality, one could:

- **Extended training and/or sampling.** Train for more epochs and/or use additional reverse diffusion steps at inference for sharper outputs.
- **Fix the reverse-process variance to an isotropic constant.** Choosing  $\Sigma_t = \beta_t I$  rather than learning a full diagonal covariance could stabilize training and improve quality.
- **Experiment with the forward noise schedule.** Tune the forward diffusion  $\{\beta_t\}$  (e.g. linear or cosine schedules) for balanced denoising.
- **Enhance U-Net architecture** Incorporate group normalization and self-attention to capture global structure and reduce artifacts.

Each of these adjustments is expected to reduce residual artifacts and improve the faithfulness of generated digits.

---

## Question 2.2 : Intermediate Reverse-Process Samples

During sampling, the DDPM’s reverse diffusion gradually transforms noise into coherent digits. At the highest timesteps ( $t = 1000$ ), each output is essentially pure noise; only around the midpoint (e.g.  $t = 500$ ) do faint contours—perhaps the rough shape of a “3” or “7”—begin to appear. It was surprising to observe that recognizable outlines emerge so late in the denoising chain. By  $t \leq 100$ , most speckles have been removed, yielding solid strokes and well-defined curves (closed loops of an “8”, continuous line of a “1”). Finally, at  $t = 0$ , the model produces a clean, fully formed digit indistinguishable from a real MNIST sample.

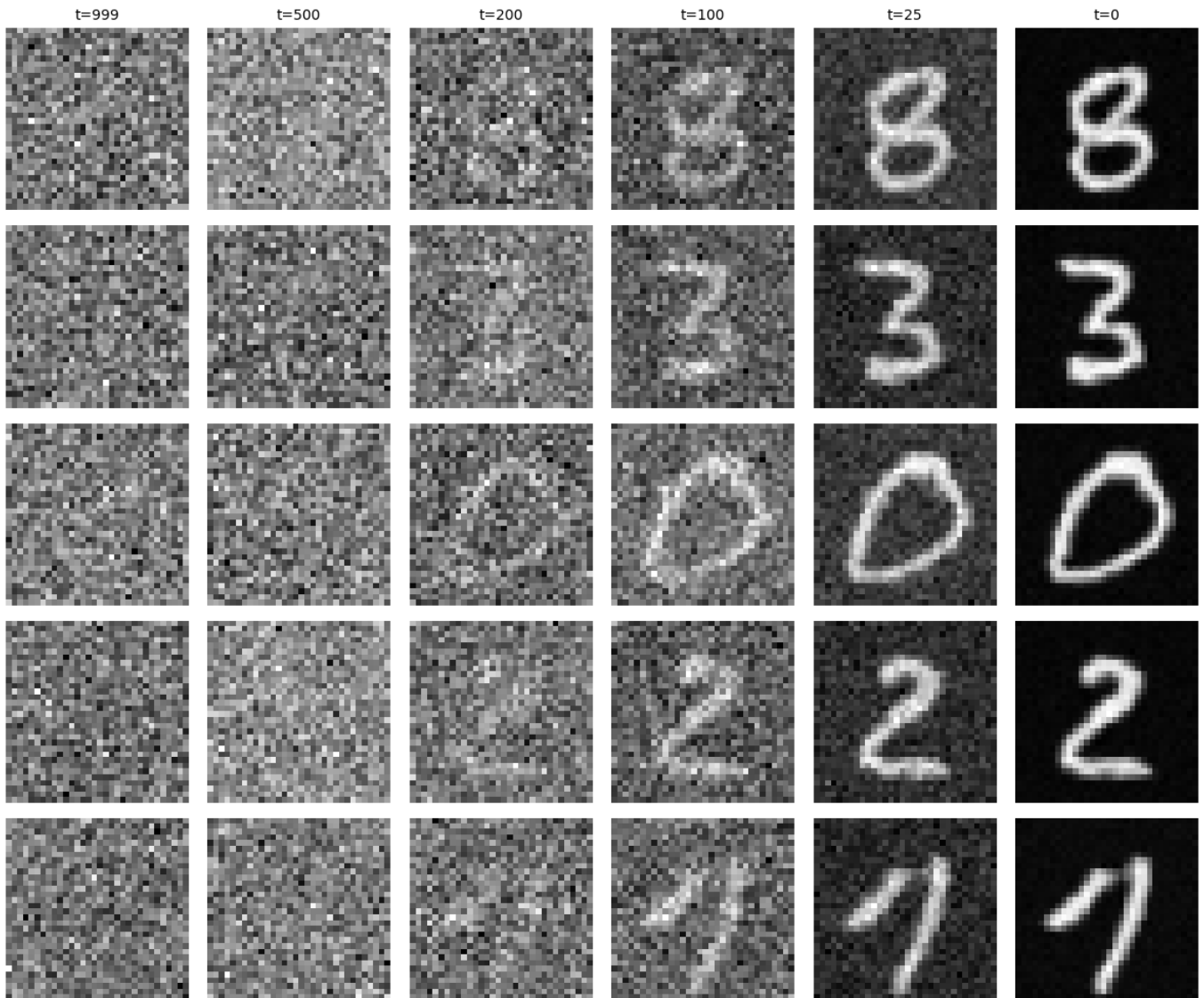


Figure 11: Intermediate reverse diffusion samples for four chains. Columns correspond to timesteps  $t = \{999, 500, 200, 100, 25, 0\}$ , and rows to different sample chains. Coherent digit shapes emerge gradually as the model removes noise.

This step-by-step refinement confirms the DDPM’s expected behavior: high-level digit structures form first, then finer details are filled in, converging to a realistic handwritten digit by the final timestep.

---

## Question 3 : Conditional Diffusion Model (CFG)

### Question 3.1: Classifier-Free and Guidance

The term “classifier-free” refers to the fact that this guidance method does not require any separately trained classifier. Instead, we jointly train the diffusion model to predict both a conditional score  $\varepsilon_\theta(z_t, c)$  and an unconditional score  $\varepsilon_\theta(z_t)$  by sometimes replacing the true label  $c$  with a null token during training. At sampling time, these two estimates are linearly combined as

$$\tilde{\varepsilon}_\theta(z_t, c) = (1 + w) \varepsilon_\theta(z_t, c) - w \varepsilon_\theta(z_t)$$

where  $w$  controls the strength of guidance. This yields the same trade-off between sample fidelity and diversity as traditional classifier gradients—but without ever training or invoking a classifier—hence “classifier-free guidance”.

### Question 3.2: Alternative to Classifier-Free Guidance

The paper points out an alternative that avoids learning an unconditional branch altogether: if the class prior  $p(c)$  is known and the number of classes is small, one can recover the marginal (unconditional) score by

$$\nabla_{z_t} \log p(z_t) = \nabla_{z_t} \log \left( \sum_c p(z_t | c) p(c) \right) = \sum_c p(c | z_t) \nabla_{z_t} \log p(z_t | c)$$

thus obtaining  $\nabla_{z_t} \log p(z_t)$  from the conditional scores alone. Under this scheme, the training loss reduces to the original DDPM denoising objective applied only to the conditional model:

$$L(\theta) = \mathbb{E} \left\| \varepsilon_\theta(z_t, c) - \epsilon \right\|_2^2, \quad z_t = \alpha_t x + \sigma_t \epsilon.$$

In contrast, classifier-free guidance uses a joint loss that alternates between conditional and unconditional updates, explicitly learning both  $\varepsilon_\theta(z_t, c)$  and  $\varepsilon_\theta(z_t)$  during training.

### Question 3.3 : Sample Generation (Across Epochs)

Introducing classifier-free guidance significantly sharpens and class-faithful digit generation as training progresses. Even in early epochs, the guided model produces noticeably clearer outputs than the VAE or unguided DDPM: although some blurring remains, each sample already bears a strong resemblance to its target class. By epochs 15–20, the digits become crisp and unambiguous, with strokes and loops well defined.

The benefits of conditioning and guidance are evident: the images are less noisy as we can see in Figure ??, looks more consistently like the intended digits showing improved fidelity to the desired digit class compared to an unguided model.

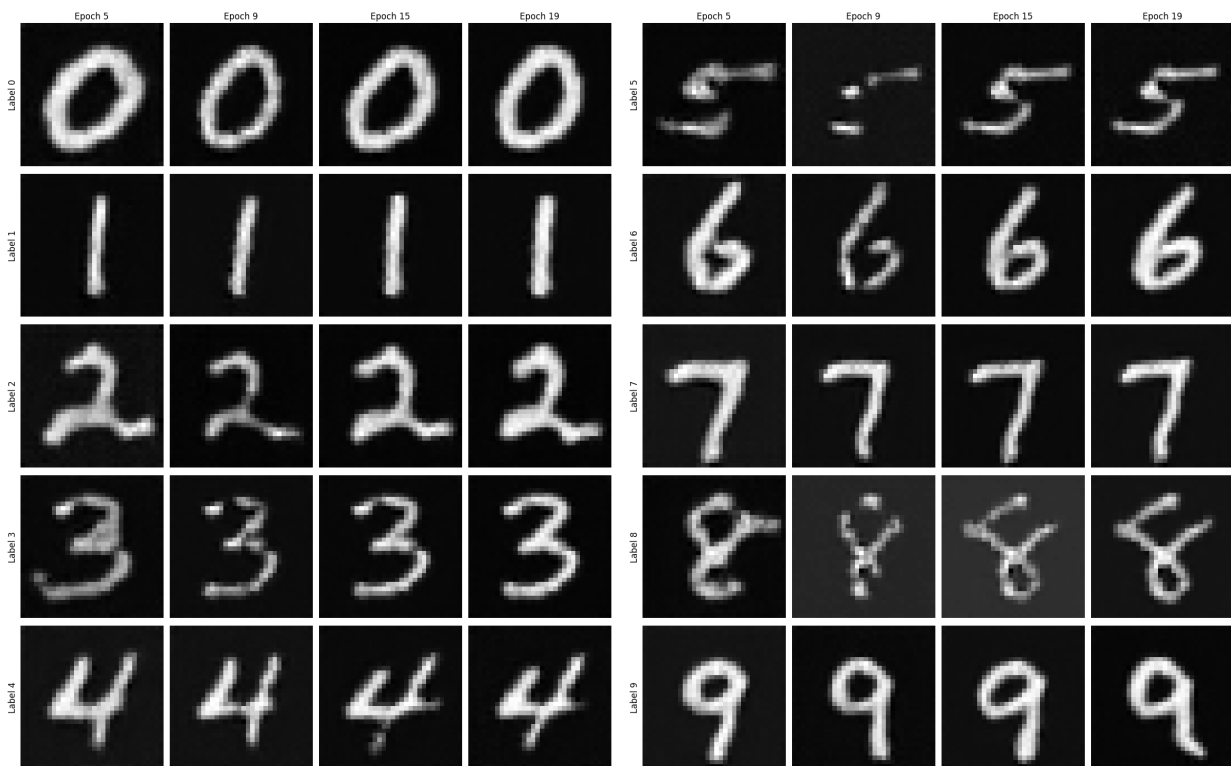


Figure 12: CFG diffusion model samples conditioned on labels 0–9 across epochs 5, 9, 15, and 19 (left panel: labels 0–4; right panel: labels 5–9). Within each panel, columns denote epochs and rows fixed digit classes, illustrating how classifier-free guidance sharpens and refines digit fidelity over training.

In summary, conditioning steers the generation process to produce well-defined digits of the requested class, and this effect becomes stronger over training: the longer the model trains, the more it can translate label information into sharper, more reliable samples for each digit category.