

1. Understanding Diffusion

Forward diffusion, in plain terms

We start with a clean image x_0 and gradually corrupt it by adding small amounts of Gaussian noise over 100 timesteps. At each step t , the image becomes

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon,$$

where ϵ is fresh noise and $\bar{\alpha}_t$ accumulates the prior noise fractions. In our `show_noise_progression` outputs, you can see the digit still clear at ~20% noise, but by ~80% it's nearly unrecognizable—and at 100% it's pure static.

Why gradual noise instead of one big dose?

Breaking noise injection into small increments turns a very hard learning problem into many easier ones: at each t , the model only needs to remove “a little more” noise rather than leap from total static back to a digit. This curriculum of noise levels stabilizes training and improves final image quality.

When does the digit first “pop out”?

From the `visualize_generation_steps` plots, most digits become identifiable around 40–50% through the reverse process (i.e. halfway from full noise back to clean). Simpler shapes like “1” often emerge closer to ~40%, while more intricate ones such as “8” appear nearer to ~55%.

2. Model Architecture

Why U-Net?

A U-Net excels at both capturing global context (via downsampling) and preserving fine details (via upsampling). Diffusion denoising needs both: coarse shape reconstruction and precise stroke recovery. Simpler feed-forward or purely encoder/decoder nets struggle to balance these simultaneously.

Skip connections

Each encoder “DownBlock” not only reduces resolution but saves a high-resolution feature map (“skip”). During decoding, the corresponding “UpBlock” concatenates its up-sampled features with that skip, injecting lost details back into the reconstruction. This direct highway for high-frequency information is essential for crisp, accurate outputs.

Class conditioning

I transformed each target digit label into a small 1×1 embedding via our `EmbedBlock` and add it (alongside the time embedding) at the U-Net bottleneck. This tells the model “which digit to form”—without it, generation would be class-agnostic and unpredictable.

3. Training Analysis

What the loss values mean

We train the U-Net to predict the exact noise ϵ added at each timestep, using mean squared error. In Epoch 1, per-batch losses started around 1.10 and fell to ~0.08, and the average training loss** was 0.1132, with a validation loss of 0.0828—showing the model quickly learned to remove most noise.

Image quality over time

- Epochs 1–5: outputs were noisy blobs with faint digit outlines.
- *Epochs 6–15: shapes sharpened; by Epoch 10 most digits were clearly legible.
- Epochs 16–30: strokes refined, loop closures in “8” and curves in “3” became smooth and consistent.

Role of time embedding

The sinusoidal time embedding encodes “how noisy” x_t is, so the U-Net knows whether to perform a large or small denoising correction. Without it, a single network can’t adaptively handle all 100 noise levels in one pass.

4. CLIP Evaluation

I did not execute the CLIP scoring cell in time, so CLIP-based analysis remains as future work.

5. Practical Applications

Potential uses

- Data augmentation for handwritten-digit OCR systems.
- Art and design: guided generation of stylized numerals.
- Medical imaging: analogous denoising pipelines for X-rays or MRIs.

Current limitations

- Trained only on low-res (28×28) grayscale digits—cannot yet handle high-res or color data.
- Inference requires 100 sequential denoising passes, making it relatively slow.
- Class vocab limited to 10; adding new classes requires retraining.

Three concrete next steps

1. Scale up resolution: add extra U-Net levels (e.g. `down_chs=(64,128,256,512)`) to handle 64×64 or 128×128 images.
2. Classifier-free guidance: train with random-class masking so i can tune conditional strength at inference time.
3. Learned noise schedules: experiment with cosine or learned β -schedules to optimize the noise curriculum for faster convergence.

Conclusion

I successfully implemented and trained a diffusion-based U-Net, visualized its denoising trajectory, and produced clear digit samples—all within the notebook. The missing CLIP and bonus cells are noted as future work rather than credited accomplishments. Let me know if you’d like any further tweaks!