

Building an Image Generator.

Recall: So far we have focused on **unconditional** generation.

Problem: Sample from p_{data}

Train: Use e.g., the conditional flow matching objective:

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{\square} \|u_t^\theta(x) - u_t^{\text{target}}(x|z)\|^2$$
$$\square = z \sim p_{\text{data}}, t \sim \text{Unif}[0, 1], x \sim p_t(x|z)$$

Sample: Simulate the corresponding ODE (or SDE):

$$dX_t = u_t^\theta(X_t)dt, \quad X_0 \sim p_{\text{init}}$$

But what about **conditional generation?**

Today's Agenda:



1. Extend our generative modeling framework from **unconditional generation** to **conditional generation**
2. Develop **classifier-free guidance** for conditional sampling
3. Discuss **architectural choices** for the prototypical case of **image generation** and **survey current models**.
4. Guest talk by Carles Domingo-Enrich!

Conditional Generation : Cond on a prompt.
Guided.

Guided Generation: What Changes?

Unguided	Guided
Marginal probability path	$p_t(x)$
Marginal vector field	$u_t^{\text{target}}(x)$
Marginal score	$\nabla \log p_t(x)$
Model	$u_t^\theta(x)$
CFM Objective	$\mathcal{L}_{\text{CFM}}(\theta)$
	Guided CFM Objective ???

Obs 1: Fix $y \Rightarrow$ Recover the unguided problem.

$$\mathcal{L}_{\text{cfm}}^{\text{guided}}(\theta) = \mathbb{E}_{\square} \left\| u_t^\theta(x|y) - u_t^{\text{target}}(x|\bar{z}) \right\|_2^2$$

$$\square = \underbrace{\bar{z}, y \sim P_{\text{data}}(\bar{z}|y)}_{\text{DataLoader}}, t \sim \text{Unif}[0,1], x \sim P_t(x|\bar{z}).$$

Obs 2: Now, let y vary.

$$\mathcal{L}_{\text{cfm}}^{\text{guided}}(\theta) = \mathbb{E}_{\square} \left\| u_t^\theta(x|y) - u_t^{\text{target}}(x|\bar{z}) \right\|_2^2$$

$$\square = \underbrace{(\bar{z}, y) \sim P_{\text{data}}(\bar{z}, y)}_{\text{DataLoader}}, t \sim \text{Unif}[0,1], x \sim P_t(x|\bar{z}).$$

Code:
DataLoader

Guided Sampling



(Prediction).

Algorithm 7 Guided Sampling Procedure

Require: A trained guided vector field $u_t^\theta(x|y)$.

- 1: Select a prompt $y \in \mathcal{Y}$, such as “a cat baking a cake”.
- 2: Initialize $X_0 \sim p_{\text{init}}$.
- 3: Simulate $dX_t = u_t^\theta(X_t|y)dt$ from $t = 0$ to $t = 1$.

Can we do better? At least empirically, the answer is yes...

Classifier-free guidance.

Recall: A Gaussian cond prob path.

$$P_t(x|\bar{z}) = \mathcal{N}(x; \beta_t^{-2} \text{Id}), \quad \alpha_1 = \beta_0 = 1, \quad \omega_0 = \beta_1 = 0.$$

FACT:

$$u_t^{\text{target}}(x|y) = \alpha_t x + \beta_t \nabla \log P_t(x|y).$$

$$(\alpha_t, \beta_t) = \left(\frac{\dot{\alpha}_t}{\alpha_t}, \frac{\dot{\alpha}_t \beta_t^2 - \dot{\beta}_t \alpha_t}{\alpha_t} \right).$$

> check previous notes.

Recall: Bayes Rule.

$$\nabla \log P_t(x|y) = \nabla \log \frac{P_t(x) P_t(y|x)}{P_t(y)} \\ = \nabla \log P_t(x) P_t(y|x) - \underbrace{\nabla \log P_t(y)}_{\parallel 0}$$

Plug in Bayes Rule,

$$U_t^{\text{target}}(x|y) = a_t x + b_t (\nabla \log P_t(x) + \nabla \log P_t(y|x))$$

$$\uparrow \quad \quad \quad U_t^{\text{target}}(x)$$

$$= \underbrace{(a_t x + b_t \nabla \log P_t(x))}_{\text{guidance with classifier.}} + b_t \nabla \log P_t(y|x).$$

$$= U_t^{\text{target}}(x) + b_t \nabla \log P_t(y|x).$$

\Downarrow
The term is a sort of "classifier".

All signals of y .

To improve the significance of y :

Choose a "guidance scale." $w > 1$.

$$U_t(x|y) = U_t^{\text{target}}(x) + \boxed{w} b_t \nabla \log P_t(y|x). \quad \begin{matrix} \rightarrow \\ \text{Plug back} \end{matrix}$$

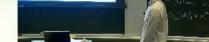
$$= U_t^{\text{target}}(x) + w b_t (-\nabla \log P_t(x) + \nabla \log P_t(x|y)) \quad \begin{matrix} \rightarrow \\ \text{Bayes Rule.} \end{matrix}$$

$$= U_t^{\text{target}}(x) - w a_t x + w a_t x + w b_t (\nabla \log P_t(x|y) - \nabla \log P_t(x))$$

$$= U_t^{\text{target}}(x) - w(a_t x + b_t \nabla \log P_t(x)) + w(a_t x + b_t \nabla \log P_t(x|y))$$

$$= (1-w) U_t^{\text{target}}(x) + w \underbrace{U_t^{\text{target}}(x|y)}$$

Classifier-Free Guidance Training



Observation: We may treat the unguided vector field as conditioned on nothing.

But, **nothing is something**:

$$u_t^{\text{target}}(x) = u_t^{\text{target}}(x|y = \emptyset)$$

We may now train a single model $u_t^\theta(x|y)$, $y \in \{\mathcal{Y}, \emptyset\}$ by re-using $\mathcal{L}_{\text{CFM}}^{\text{guided}}(\theta)$ and occasionally setting $y = \emptyset$:

$$\mathcal{L}_{\text{CFM}}^{\text{CFG}}(\theta) = \mathbb{E}_\square \|u_t^\theta(x|y) - u_t^{\text{target}}(x|z)\|^2$$

$\square = (z, y) \sim p_{\text{data}}(z, y)$, with prob. η , $y \leftarrow \emptyset$, $t \sim \text{Unif}[0, 1]$, $x \sim p_t(x|z)$

只需 train 一个 model.

在 train 过程中随机插入 non-prompt.

$\Pr[y = \emptyset] = \eta$. Randomly insert.

Guidance: $w = 1.0$
0 0 0 0 0 4 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 5 2 2 2 2
3 2 1 3 3 3 1 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 1 5 5 8 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 1 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 1 9 4 9 9 9 9 9
0 1 3 2 9 8 4 0 2 7

Guidance: $w = 3.0$
0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9
2 1 8 9 2 0 5 5 9 6

Guidance: $w = 5.0$
0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9
1 2 3 4 7 1 5 0 2



three!

Algorithm 8 Classifier-Free Guidance Sampling Procedure

Require: A trained guided vector field $u_t^\theta(x|y)$.

1: Select a prompt $y \in \mathcal{Y}$, or take $y = \emptyset$ for unguided sampling.

2: Select a **guidance scale** $w > 1$.

3: Initialize $X_0 \sim p_{\text{init}}$.

4: Simulate $dX_t = [(1 - w)u_t^\theta(X_t|\emptyset) + wu_t^\theta(X_t|y)] dt$ from $t = 0$ to $t = 1$.

$w \uparrow \Rightarrow$ clear generation.
 tradeoff between { quality diversity } of generation.

Architecture. (for $u_t^\theta(x|y)$)

Architectures for Image Generation



Recall: An image lives in $\mathbb{R}^{C_{\text{image}} \times H \times W}$

Question: An MLP is insufficient in such a high-dimensional space.
What, then, should $u_t^\theta(x|y)$ look like?

Preview: We'll explore two choices: **U-Nets** (convolution based) and **diffusion transformers** (attention based).

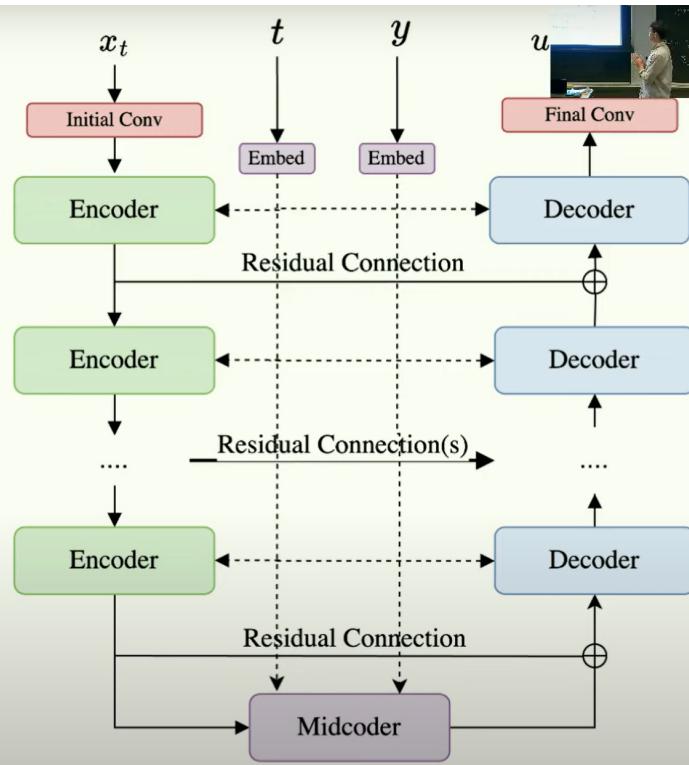
Pay Attention: How is y encoded, embedded, and processed?

Lab Three U-Net

In lab three, we'll utilize the simplified **U-Net architecture** shown at right to build a generative model for the **MNIST dataset**.

In this case $x_t \in \mathbb{R}^{1 \times 32 \times 32}$ and

$y \in \{0, 1, \dots, 9, \emptyset\}$



Diffusion Transformer (DiT)

paper [1] (right).

Idea: Divide an image into **patches** and **attend** between the patches. Based on the **vision transformer** (ViT).

