

Diffusion

- decomposing the mapping into a sequence of small, incremental transformations
 - from a simple distribution (e.g., Gaussian noise)
 - to a complex data distribution
- learn to reverse a gradual noising process
 - instead of directly learning a single complex mapping from noise to data

Denoising Diffusion Probabilistic Models (DDPM)

- probabilistic diffusion model that defines a Markovian forward noising process
 - data distribution → gradually corrupted by Gaussian noise → normal Gaussian
 - close-form transition in the noising (forward) process
- learns to reverse the noising process step by step using a neural network
 - gradually denoising Gaussian noise estimated by the network (denoiser)

$$q(x_t | x_0) = \mathcal{N}(\sqrt{\alpha_t} x_0, (1 - \alpha_t) \mathbf{I})$$

linear Gaussian Markov chain
Gaussian

$$q(x_{t-1} | x_t, x_0) \approx p_\theta(x_{t-1} | x_t)$$

Gaussian (tractable) ← *approximate* *infinitesimal change → Gaussian noise*

Denoising Diffusion Implicit Models (DDIM)

- assume a non-Markovian forward process
- defines a deterministic mapping from x_t to x_{t-1}
 - prediction → decomposed into clean data x_0 + noise vector to $t - 1$ (+ stochastic term)

- skipping steps during sampling → faster inference with fewer steps (accelerating)

Classifier-Free Guidance (CFG)

$$\epsilon = \epsilon_{\text{uncond}} + s \cdot (\epsilon_{\text{cond}} - \epsilon_{\text{uncond}})$$

- conditioning technique for diffusion models (conditional generation)
 - not require an external classifier
- by Bayesian theorem
 - $p(x_t | y) \propto p(y | x_t) p(x_t)$
 - $\nabla_{x_t} \log p(x_t | y) = \nabla_{x_t} \log p(x_t) + \nabla_{x_t} \log p(y | x_t)$
 - interpreted as approximating the Bayesian posterior ($\nabla_{x_t} \log p(x_t | y)$)
 - by amplifying the likelihood term ($\nabla_{x_t} \log p(y | x_t)$) without an explicit
- how it works
 - train the model with and without conditioning
 - during sampling → combine conditional and unconditional predictions
 - guidance scale $s \uparrow$ → stronger conditioning
 - too large scale \uparrow → diversity \downarrow
- pushing samples toward regions that better satisfy the condition
 - interpolating between unconditional and conditional noise predictions

Score Function

- $\nabla_x \log p(x)$: gradient of log density
 - how to move a noisy sample toward higher-density regions of the data distribution.

Tweedie's formula

$$\mathbb{E}[x_0 | x] = x + \sigma^2 \nabla_x \log p(x)$$

- posterior mean estimation using variance-scaled score function
 - denoising \approx estimating the posterior mean using the score of the noisy distribution
 - in the diffusion models \rightarrow how far I should move in that direction to recover the clean signal

Flow Matching

- reframes generative modeling as learning a continuous transport
 - from a simple distribution to the data distribution *push forward*
 - velocity regression without explicitly optimizing likelihoods or solving expensive ODEs during training
- variants dependent on how to define a conditional flow mapping (transformation)
 - diffusion models: a special case of flow matching
 - where the target velocity corresponds to the score-induced reverse-time dynamics

Normalizing Flow

composition of invertible transforms

- learns a single invertible mapping *transformation Φ 이 invertible 해야 함*
 - from a simple distribution (e.g., Gaussian) to the data distribution *Jacobian의 determinant 계산이 어려움 $O(n^3)$*

- uses the change-of-variables formula: *산각 하에서 elementwise 곱 구함 (\therefore complexity \downarrow)*

$$\log p_1(x_1) = \log p_0(\Phi^{-1}(x_1)) + \log |\det(J)|$$

◦ the transformation Φ must be invertible

◦ the Jacobian J determinant must be tractable

- limitation: computing Jacobians is expensive

$$J = \frac{\partial \Phi^{-1}(x_1)}{\partial x_1}$$

*locally linear approximation
determinant \rightarrow 행렬의 scaling*

확률 밀도 보존

$$p_X(x) dx = p_Z(z) dz$$

$$p_X(x) = p_Z(z) \left| \det \left(\frac{dz}{dx} \right) \right|$$

Coupling layers: 여러 층의 transformation 쌓아 적음

NLL 최소화 (log likelihood 최대화)

부피 변화율에 의해 밀도 재조정

Continuous Normalizing Flow

ODE + log density

- extends normalizing flows to continuous time
 - as the integration of infinitesimal changes instantaneous
- models the transformation as an ODE model이 학습하는 건 변환 그 자체가 아니라 ODE 학습
 - $dx/dt = u_t(x)$ timestep t에서 x 지점에서 어디로 이동할지
 - density evolution follows the continuity equation:
 - $\partial p_t(x)/\partial t = -\nabla \cdot (p_t(x)u_t(x))$ timestep t에서 x 지점에서 확률밀도 p가 어떻게 변할지
 - log-density change: trace: 각 축 방향으로의 infinitesimal scaling ≡ 연산
 - $d \log p(x_t)/dt = -\nabla \cdot u_t(x_t) = -\text{Tr} \left(\frac{\partial u_t}{\partial x} \right)$ divergence를 선에 대해 적분 integral 요구됨
- advantages
 - more flexible transformations
 - avoids explicit Jacobian determinants $O(d^2) \rightarrow$ 근사하여 확률적 추정
- drawbacks
 - training requires solving ODEs repeatedly
 - computationally expensive

$$x_{t+\Delta t} \approx x_t + u_t(x_t, t) \Delta t$$

log likelihood 최대화로 풀거나
vector regression으로 풀거나

Flow Matching

- matching vector fields instead of densities
 - avoids likelihood-based training altogether
- learns a neural velocity field $v_\theta(x, t)$ to match a target velocity $u_t(x)$ transport 우리가 path를 어떻게 설정하느냐에 따라 target velocity 구할 수 있음
- objective
 - $\mathcal{L} = E[||v_\theta(x_t, t) - u_t(x_t)||^2]$
- advantages
 - no need to compute log-densities
 - no need to solve ODEs during training
- drawbacks
 - the marginal velocity field $u_t(x)$ is generally intractable GT velocity field

Conditional Flow Matching

- learn a conditional velocity field $u_t(x|x_0, x_1)$
 - instead of learning a marginal velocity field $u_t(x)$
- samples are paired:
 - $x_0 \sim$ base distribution
 - $x_1 \sim$ data distribution
- learns how to transport x_0 to x_1 over time
- objective
 - $\mathcal{L} = E[||v_\theta(x_t, t) - u_t(x_t|x_0, x_1)||^2]$
 - theoretically grounded by the continuity equation
- benefits
 - simpler target vector field
 - stable training
 - avoids density estimation entirely

Rectified Flow

transport path를 어떻게 정의하는지에 따라 velocity가 달라짐

- special case of conditional flow matching
- defines straight-line paths between x_0 and x_1 :
 - $x_t = (1 - t)x_0 + tx_1$

e.g. diffusion-style

$$X_t = \alpha(t) X_1 + \sigma(t) X_0$$

→ 얼마나 큰 sign을 보존할지

↳ noise를 얼마나 추가할지

$$u_t(X_t|X_1) = \frac{dX_t}{dt} = X_1 - X_0$$

- corresponding velocity field is simple and deterministic
- advantages
 - extremely simple training
 - fast sampling
 - strong empirical performance
- widely used in practice due to its simplicity

L2 비용 함수 (cost function)의 OT에서는

각 점이 직선으로 이동하는 geodesic에 가까움

↪ 직선 → optimal transport에 가까울수록
coupling에 따라 비용 달라짐

OT-optimal coupling이어야 직선이 OT

RD에서는 OT-optimal coupling 찾는 건 어려지만

더 일관된 transport → 더 OT에 가까운 성질 유도

Diffusion Transformer (DiT)

- latent diffusion VAE from SD \because pixel DDPM \rightarrow computationally expensive
- patchify \rightarrow token sequence patch size $\downarrow \rightarrow$ # of tokens $\uparrow \rightarrow$ computational cost: memory \uparrow
quality \uparrow
- positional embedding (sin/cos)

- conditioning timestep t condition c

in-context concatenate t, c into input token sequences

cross-attention

adaLN layer norm의 scaling \cdot shift \leftarrow regression from condition

adaLN-zero residual block — initially zero (gating) training stability \uparrow

G-Flops * of parameters 보다 지표로서 좋음 (good proxy)

forward pass \because attention은 token 수에 민감 (이러지 핵심도에 따라 요구되는 연산 달라질 수 있음)

SDE-based Score Model

forward process 정의: data \rightarrow noise

$$dx = \underbrace{f(x,t)}_{\text{drift}} dt + \underbrace{g(t)}_{\substack{\text{diffusion} \\ \text{coefficient}}} \underbrace{dW_t}_{\substack{\text{diffusion} \\ \text{(stochastic term)} \\ \text{Brownian motion}}} \quad dW_t \sim N(0, I)$$

reverse process

reverse-time SDE

$$dx = \underbrace{[f(x,t) - g(t)^2 S_t(x)]}_{\text{drift}} dt + g(t) d\bar{W}_t \quad S_t(x) = \nabla_x \log P_t(x)$$

score function

probability flow ODE

$$\frac{dx}{dt} = F(x,t) = f(x,t) - \frac{1}{2} g(t)^2 S_t(x) \quad \text{deterministic}$$

"diffusion sampling은 ODE/SDE를 적분하는 문제"

VP-SDE (Variance - Preserving)

forward 과정에서 x_t 의 variance가 폭발하지 않도록

$$SNR(t) = \frac{\text{Var}(\text{signal})}{\text{Var}(\text{noise})} = \frac{\alpha_t}{1 - \alpha_t}$$

forward SDE: $dx = -\frac{1}{2} \beta(t) x dt + \sqrt{\beta(t)} dW_t$

drift: $f(x,t) = -\frac{1}{2} \beta(t) x$

시간에 따라 signal 감쇠

diffusion: $g(x,t) = \sqrt{\beta(t)}$

DDPM의 discrete-time forward \rightarrow continuous version

$$x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon$$

$$\text{Var}(x_t) = \alpha_t \text{Var}(x_0) + (1 - \alpha_t) \text{Var}(\epsilon) = 1$$

t 증가에 따라 signal의 분산 \downarrow

timestep별 SNR 편차 심함 \rightarrow 학습 신호 불균형

낮은 SNR 구간에서 예측 어려움

reverse 시

PF-ODE: $\frac{dx}{dt} = -\frac{1}{2} \beta(t) x - \frac{1}{2} \beta(t) S_0(x,t)$

noise로부터 유도 가능

$$S_0(x,t) \propto -\frac{1}{\alpha(t)} \epsilon_0(x,t)$$

VE-SDE (Variance - Exploding)

forward 과정에서 x_t 의 variance \uparrow

$$SNR(t) \propto \frac{1}{\sigma(t)^2}$$

forward: $dx = g(t) dW_t$

drift: $f(x,t) = 0$ diffusion: $g(x,t) = \sqrt{\frac{d}{dt} \sigma(t)^2}$

NCSN

$$x_t = x_0 + \sigma(t) \epsilon$$

perturbed data

여러 각 σ_t 에서 score 학습

annealed Langevin Dynamics

X - Prediction

$$x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \varepsilon \quad \varepsilon \sim N(0, I)$$

* ε -prediction

score matching과 정합 (noise \equiv score function의 scaled 역방향)

noise $\varepsilon_\theta(x_t, t)$ 자체를 예측

timestep 따라 loss scale 달라질 수 있음 (\therefore SNR)

* x_0 -prediction

clean data \hat{x}_0 자체를 예측

최종 objective와 정합

작은 SNR 구간에서는 x_0 추정 어려움 \rightarrow 불안정

* V-prediction

ε 와 x_0 의 선형 결합

두 신호 섞어서 SNR 균형 맞춤

$$v = \sqrt{\alpha_t} \varepsilon - \sqrt{1 - \alpha_t} x_0$$

Noise Scheduler

forward process의 각 t에서 $\alpha_t, \bar{\alpha}_t$ 정합 \rightarrow noise 양 결정

Coefficient

$$\text{SNR} = \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t} \quad \bar{\alpha}_t \downarrow \rightarrow \text{noise} \uparrow \rightarrow \text{SNR} \downarrow \rightarrow \varepsilon \text{ 예측 어려움}$$

loss를 단순 MSE로 두면 unstable

cosine schedule

SNR이 너무 빨리 붕괴하는 구간 막음

Sampler와 함께 고려되어야 함 (\therefore 어떤 구간에서는 abruptly rapidly 변화 가능)

오차가 민감한 구간에 step 더 배치 (SNR를 고려해서)

Sampler

Update 규칙

DDPM forward process : Markovian chain + linear Gaussian transitions + $q(x_t|x_0)$ closed form

reverse process : stochastic reverse \rightarrow noise 추가 \nearrow 각 step에서 분포에서 직접 뽑는 형태
ancestral sampling $p(x_{t-1}|x_t) = N(\mu(x_t, t), \sigma_t^2 I)$ (Markovian chain)

이산적으로 sampling 하는 게 핵심

DDIM forward process : forward process 자체는 marginal에 대해 DDPM과 동일 (실무에서는 transition 자체도 DDPM처럼)

reverse process : implicit reverse process가 Markovian chain이 아님 (x_0 로 가는 게 고려되므로 joint distribution)

implicit trajectory (mapping) x_0 예측 + σ_t 로 가는 방향 예측 with scaling + stochastic term
probability flow ODE로 표현 가능 (ODE sampling) 0이면 deterministic
deterministic 경로 solver step $\downarrow \rightarrow$ accelerate \uparrow but error \uparrow

DDPM과 같은 marginals $q(x_t|x_0) = N(\sqrt{\alpha_t}x_0, (1-\alpha_t)I)$ 를 유지하는 diffusion family 만들

이를 만들기 위해 posterior $q_{\theta}(x_{t-1}|x_t, x_0)$ 분포를 정의 (여러 개가 있지만 그중 하나)

DDPM의 forward처럼 쓰면 결국 있는 유도되는 거 아닌가? \rightarrow 핵심은 marginal $q(x_t|x_0)$ 를 맞추는 것이지
DDPM처럼 $q(x_t|x_{t-1})$ 를 Markov transition으로 고정하는 게 X

(실무에서 forward 자체는 DDPM처럼 but reverse process에서는
ancestral Markov chain 라는 다른 implicit trajectory 사용)

요약) DDIM은 DDPM과 동일한 forward marginals $q(x_t|x_0)$ 를 유지하는 diffusion family 고려
그 family 안에서 경로의 stochasticity를 고려하는 파라미터 σ_t 를 두고, 그에 맞는 $q_{\theta}(x_{t-1}|x_t, x_0)$ 를 정의
정의를 전개하면 sampling update가 deterministic term + stochastic term

Solver

ODE/SDE를 step-by-step으로 근사해 푸는 수치 해석

Sampler의 dynamics를 기반으로 얼마나 정확할지 (주로 DDIM/ODE/SDE 기반의 가역 sampling 시 사용)

deterministic \rightarrow probability flow ODE로 고려 가능

model의 예측값으로부터 drift 포함 \rightarrow sampler 규칙에 수치적으로 근사

1st Order Solver (Euler 계열)

한 step에서 현재 시점의 기울기 (derivative)만 보고 직선으로 한번 이동 (denoising process)

ODE $\frac{dx}{dt} = F(x, t)$ 일때 Euler $x_{k+1} = x_k + h F(x_k, t_k)$
denoised \quad step size h \approx timestep 차이

e.g. DDIM PF-ODE 관련 reverse process 시

$$\frac{dx}{dt} = F(x, t) = -\frac{1}{2} \beta(t) x - \frac{1}{2} \beta(t) S_0(x, t)$$

$S_0(x, t)$ 예측 \rightarrow score $S_0(x, t)$ 구함 \rightarrow drift $F(x, t)$ 계산 \rightarrow 다음 state update

SDE: Euler - Maruyama

$$dx = f(x, t) dt + g(t) dW_t$$

$$x_{k+1} = x_k + h f(x_k, t_k) + g(t_k) \sqrt{h} z \quad z \sim N(0, I)$$

Higher-Order Solver 한 step에서 한 기울기만 보는 게 아니라 중간 or 다음 시점도 한번 더 보도록 보정

e.g. $x_{k+1} = x_k + \frac{h}{2} (F(x_k, t_k) + F(\tilde{x}_k, t_{k+1}))$

step 수 v.s. NFE (Number of Function Evaluations)

N차 solver \rightarrow 고풍 더 늘어남

PLMS (Pseudo Linear Multistep)

DDIM update 기반 이전 step들의 ϵ 등 예측 저장 \rightarrow 다음 update 정확히

DPM - Solver

diffusion의 ODE를 단순 Euler로 풀지 말고, 해당 ODE 구조에 맞춘 고차 근사식 사용

10 ~ 20 step \rightarrow well

CFG scale $\uparrow \rightarrow$ drift $\uparrow \rightarrow$ step 배치 촘촘하게 or 안정적인 solver

ADM (Diffusion Models Beat GANs)

Classifier Guidance classifier의 $\nabla_x \log p_\theta(y|x)$ 를 이용

reverse variance 학습

$$p_\theta(x_{t+1} | x_t) = (\mu_\theta(x_t, t), \sigma_\theta^2(x_t, t) I)$$

범위 clip하여 제한해서 학습 $[\log \tilde{\beta}_t, \log \beta_t]$

→ 더 작은 step으로도 품질 유지 가능

Diffusion Likelihood

많은 step의 ϵ transition의 곱 → exact likelihood 계산 어려움

< discrete ELBO
continuous score-SDE/probability flow ODE CNF 처럼 likelihood 계산

Efficient DiT

