

## Paper summary

AI VISION Lab

1. 공부한 논문의 제목, 게재된 학회 혹은 저널 등 논문 기본 정보를 적으세요.
  - A. 이름: EDM: Elucidating the Design Space of Diffusion-Based Generative Models
  - B. 저널: NeurIPS
  - C. 도메인: Diffusion
  - D. 출판연도: 2022
  - E. 저자: NVIDIA (Tero Karras et al)
2. 논문에서 제안한 알고리즘 및 프레임워크에 대해 본인이 이해한대로 다이어그램을 그려보세요. 논문 Figure를 그대로 따라 그리면 안됩니다.
  - A. [선행연구] SDE, ODE를 제안한 Yang song의 논문을 바탕으로 전개함.
  - B. 본 저자는, 기존의 Diffusion Model (VE, VP)의 Framework의 경우는 서로 의존성이 존재하다 보니, Model의 전체구조를 수정하는 것이 아니라면 변화를 주기 어렵다고 판단함.
  - C. 따라서, Diffusion Model의 구성요소들에 대해서, 발전된 구조와 함께, 각 요소 사이에 존재하던 의존성을 제거하고자 함.
  - D. SDE, ODE의 경우 1 Step의 변화를  $dx$ 로 설명하고 있음. 그러나 해당  $dx$ 의 경우,  $f(x, t)$  - drift항과,  $g(t)$  - diffusion항으로 설명하고 있으며 이는 DDPM 기준,  $\beta$ 에 대한 식으로 설명 가능함.
  - E. 그러나 실제 Diffusion을 적용하는 과정에서는,  $dx$ 가 아닌,  $p_{0t}$ 를 계산해서 사용하고 있다는 점에서, 계산이 여러 번 필요하며 설명이 복잡하다는 한계가 있었음
  - F. 따라서,  $dx$ 에 대한 식에서 출발하는 것이 아닌,  $p_{0t}$ 를  $s(t)$ 와,  $\sigma(t)$ 로 정의한 후,  $dx$ 를 정의함. 즉, Diffusion의 모든 요소들을  $s(t)$ 와,  $\sigma(t)$ 로만 설명할 수 있게 됨.

- G. ODE의 경우, Stochasticity 가 없는 구조다 보니, 정확하고, Sampling을 효과적으로 할 수 있다는 장점이 있음. 따라서, 논문에서는 우선 ODE 기반의 Sampler를 제안함.
- H. 기존의 Sampler로 가장 많이 사용되는 ODE Solver는 Euler's Method이며, 이는 1<sup>st</sup> 에 해당함. 즉, Trajectory를 예측해야 하는 Solver의 특성상, 예측 정확도(Truncation Error)가 낮을 수밖에 없음. 또한, Runge-Kutta's method의 경우 high order의 Solver로, 예측 정확도가 높으나, 연산요구량이 상당히 많다는 한계가 있었음. 따라서, 논문에서는 정확도를 높이되, Step에 정확히 비례하는(추가 계산이 거의 없는) Trade-off(quality – NFE:  $D_\theta$ 를 몇 번 호출했는가?) 가 가능한 Sampler를 제안하고자 함. 따라서, 2<sup>nd</sup> Heun's method를 제안함.
- I. SDE의 경우, ODE 와 달리 Stochastic하다는 점 때문에 다루기 어렵다는 한계가 있으나, ODE 대비 성능이 높다는 장점이 있고, 논문에서는 ODE에서 제시한 구조를 기반으로 SDE Sampler를 정의함. (2<sup>nd</sup> Heun + Langevin Like)
- J. 또한, Training과정에 대해서도 정의하는데, 기존의 Model은 Sigma Level에 영향을 직접적으로 받는 구조라는 한계가 있었음. 또한, 이 논문에서는, NN이 직접  $D_\theta$ (Score Function of P를 기반으로 Noise를 제거한 Data)를 예측하는 것이 아닌, Score Function of P를 예측하고, 이를 기반으로 D를 구하는 구조를 제안함. (안정적)
- K. 결과적으로, 제안한 모든 Option을 결합했을 때, 가장 높은 성능을 보임.
3. 본인이 생각하는 이 논문의 장점이 무엇이라고 생각하나요? **논문 Contribution bullet을 그대로 따라 적으면 안됩니다.**
- A. Noise Scheduler, Mean, Variance 구조 간의 의존성에서 벗어나, 각 구성 요소를 제어할 수 있는 Framework를 제안했다는 점과, 기존의 Euler's method가 아닌, Heun's method를 제안함으로써, Sampling 성능을 향상시키거나, dx를 통합한 형태에서 정의함으로써, Parameter에 대한 통일성을 부여했다는 점 등, 결과적으로 Diffusion Model의 다양한 구성요소들에 대해서, 하나씩 기여했다는 점이 가장 큰 장점이라고 생각함.
4. 이 논문을 읽으면서 느낀 점, 혹은 배운 점이 있으면 적어보세요.

- A. 선행연구와 비교과정을 상세히(식, 그래프) 설명해 둬. 이를 바탕으로 추후, 복습할 일이 있는 경우 논문을 다시 읽어보면 좋을 것으로 판단됨.
5. 이 논문의 한계점이 있다면 무엇이라고 생각하나요?
- A. Tuning할 hyperparameter의 수가 상당히 많은 반면, 논문에서 제안한 Optimal한 구조의 경우, ImageNet, CIFAR10 등의 데이터셋을 기준으로 성능을 뽑았다는 한계가 있음. 즉, 다른 Task에 대해서는 이 논문과 다른 Hyperparameter를 찾아야 한다는 한계가 있음.
6. 본인의 연구에 접목시켜볼 점이 있을지 생각하고 적어보세요.
- A. Diffusion 선행연구
7. 본 Summary를 작성하는 과정에서 생성형AI를 사용했나요?
- A. 아니요

날짜: 2025-07-13

이름: 신준원

# Elucidating the Design Space of Diffusion-Based Generative models 01/13

## Introduction

이론에 기반한 diffusion model



Design space 상, 특정 범위에 어려움으로 작용  
본 논문 허용 한계는 해결할 수 있는 방법을 제안함.

- ① 이론보다는 확률 측면에서, 다양한 구성 (forward, reverse, scheduling...)  
↳ system을 부차적. option 변경 불가능한 구조.

⇒ 각 구성안의 이론적 파악 + 이론적 상제인 경우 교환 시도.

- ② 내부적인 sampling 속도.

⇒ 가장 성능 좋은 sampler는? → sampling 속도의 비약.

- ③ Score 기반 modeling

①, ②, ③ 은 분석 ⇒ Diffusion 성능 향상에 기여.

## notation

$P_{\text{data}}(x)$  : data distribution

$\sigma_{\text{data}}$  : standard deviation → noise의  $\langle \epsilon \rangle$ .

$P(x; \sigma)$  : mollified distribution → 추론 분포 (forward) :  $P(x; \sigma_{\text{max}}) \sim N(0, I \sigma^2)$   
(forward →  $\sigma_0 < \sigma_1 < \dots < \sigma_{\text{max}} \sim N(0, \sigma_{\text{max}}^2 I)$ )

$$P_{\theta_t}(x(t) | x(0)) = N(x(t); S(t)x(0), S(t)^T \sigma(t)^2 I)$$

$$S(t) = \exp\left(\int_0^t f(\epsilon) d\epsilon\right), \quad \sigma(t) = \sqrt{\int_0^t \frac{\theta(\epsilon)^2}{S(\epsilon)^2} d\epsilon}$$

$$\rightarrow x = S(t) \cdot x_0 + \sigma(t)$$

# ODE formulation

①  $f(x, t) = f(x) \in \mathbb{R}^d$ , ②  $d = d_t$ .

original ODE (probability flow ODE) → by Yang-gong 2021

$$dx = \left[ f(x) - \frac{1}{2} f(x)^2 \nabla_x \log p_t(x) \right] dt.$$

→ 실제 계산  $\Rightarrow dx(x)$  marginal distribution

$\Rightarrow$  그렇다면  $d$ 에 정의가 아님,  $S(t), \beta(t)$ 를 정의 (marginal  $p_t(x)$ )  $\Rightarrow$  이 흐름!

$$p_t(x) = \int_{\mathbb{R}^d} p_{data}(x_0) p_{data}(x_0) dx_0$$

$$= \int_{\mathbb{R}^d} p_{data}(x_0) \cdot \left[ \mathcal{N}(x; S(t)x_0, S(t)^2 \beta(t) I) \right] dx_0$$

$$= \int_{\mathbb{R}^d} p_{data}(x_0) \left[ S(t)^{-d} \mathcal{N}(x/S(t); x_0, \beta(t) I) \right] dx_0$$

$$= S(t)^{-d} \int_{\mathbb{R}^d} p_{data}(x_0) \mathcal{N}(x/S(t); x_0, \beta(t) I) dx_0$$

$$\rightarrow u = S(t) \cdot x_0, \quad x_0 = \frac{u}{S(t)}, \quad dx_0 = S(t)^{-d} du$$

$$= S(t)^{-d} [p_{data} * \mathcal{N}(0, \beta(t) I)](x/S(t))$$

$$\Rightarrow p(x; \beta) = p_{data} * \mathcal{N}(0, \beta(t) I) / p_t(x) = S(t)^{-d} p(x/S(t); \beta(t))$$

$$\Rightarrow dx = -\dot{\beta}(t) \beta(t) \nabla_x \log p(x; \beta(t)) dt$$



### Denoising Score Matching

$$E_{y \sim \text{data}} E_{n \sim N(0, \sigma^2 I)} [\|D(y+n; \theta) - y\|_2^2], \quad \nabla_{\theta} \log P(x; \theta) = (D(x; \theta) - x) / \sigma^2$$

$\rightarrow$  noise step  $\times$  data  $\rightarrow$  loss  $\Rightarrow D_{\theta}(x; \theta)$

①  $y = \text{data (image)}$

②  $n = \text{noise}$

③  $D(x; \theta) = \text{denoiser function}$

### Time-dependent signal scaling

$$x = \frac{s(t)}{s(t)} \cdot \hat{x} \quad \left. \begin{array}{l} \downarrow \\ \text{Scale scheduler} \end{array} \right\} \rightarrow \text{PDF, ODE solution. 등}$$

ODE - solution 등

$$dx = \left[ \frac{\dot{s}(t)}{s(t)} x - s(t)^2 \cdot \dot{s}(t) \cdot \nabla_x \log P\left(\frac{x}{s(t)}; \theta(t)\right) \right] dt$$

~~Time-dependent signal scaling~~

Solution by discretization

이러한  $\nabla_x \log p \rightarrow$  ODE - solution 등

④ 2<sup>nd</sup> order solver

## deterministic sampling (improved)

sample quality 상승 = sampling 위한 영향 감소.

가설 1) Sampling 과정은 독립적이다.

→ 실제 정의한  $S_0$  은  $S(0)$ ,  $S(1)$  등과 의존성 X.

model (score) 즉, 서로 black box 한지

실험) 3개의 모델 (DDPM++ : VP, NCSN++ : UZ, ADM)

① sampler 적용 (original or paper).

→ 의존성이 없다면, 성능↓

의존성이 있다면, 3 모델 모두에서 성능향상

결과) 3 모델의 sample 성능 상승, ⇒ 독립

② VFZ 비교 : 샘플 몇번 호출, ⇒ Sampling 횟수를 의미 (비례)

본문 ⇒ improved version 적용시,

3 model 모두에서 성능향상 확인 가능

→ VFZ가 적은 경우에도 유의미한 FID 관찰 가능







## Trajectory curvature & noise schedule

현재  $\Rightarrow$  ODE 정의 with  $b(t), s(t)$

$\hookrightarrow$  크기 예측 가능 ( $dy/dt$ )

Best option  $\Rightarrow b(t) = \pm 1$  &  $s(t) = 1$

$$\frac{dy}{dt} = (1 - b(a; t)) / t \rightarrow a, t \text{ fixed} \rightarrow \text{noise}$$

Sampling (with Euler step  $\rightarrow t=0$ )

$VZ$

$Vp$

$$b(t) = \pm 1, s(t) = 1$$

① trajectory 곡률 제어

$\Rightarrow$  noise 가짐. (작은 차이  $\rightarrow$  곡률)

$\frac{2}{7}$ , 최적화한 개체 유지 = Best

$\hookrightarrow$  linear 한 이동 = 성능 우수

## Stochastic sampling

ODE sampling  $\rightarrow$  forward-reverse sampling.

또한 SDE 24비 샘플이  $\frac{1}{\lambda_0}$

Q) ODE, SDE 모두 같은 distribution을 따르는데, 왜 굳이 stochasticity를 부여하는가?

$$dx_t = \underbrace{-\dot{\beta}(t) \beta(t) \nabla_x \log p(x; \beta(t)) dt}_{\text{Probability flow ODE}} + \underbrace{\pm \beta(t) \nabla_x \log p(x; \beta(t)) dt + \sqrt{2\beta(t)} dW_t}_{\text{Langevin diffusion SDE}}$$

## denotation

$dx_+ \Rightarrow$  forward  
 $dx_- \Rightarrow$  reverse

Langevin term.

$\Rightarrow$  deterministic score-based denoising term (reverse)

+ stochastic noise term.

$$\beta(t) = \frac{\dot{\beta}(t)}{\beta(t)} \Rightarrow \text{제곱의 } \beta(t) \Rightarrow \frac{1}{\lambda_0}$$

## Paper $\rightarrow$ stochastic sampler

Alg:  $2^{nd}$  order deterministic ODE integrator + Langevin like

$$\text{Alg 1) } r_i = \begin{cases} \min(\frac{Sch_{\text{min}}}{N}, \sqrt{\epsilon} - 1) & \text{if } t_i \in [t_{\text{min}}, t_{\text{max}}] \\ 0 & \text{otherwise} \end{cases}$$

$$\hat{t}_i = t_i + r_i t_i \rightarrow \hat{x}_i \leftarrow x_i + \sqrt{\hat{t}_i - t_i} z$$

Alg 2)  $\hat{x}_i, \hat{t}_i$  기반  $\rightarrow$  dd step (single step)

$$\text{Euler step } x_{i+1} = \hat{x}_i + (\hat{t}_{i+1} - \hat{t}_i) d_i$$

$\rightarrow$  order correction

return  $x_N$



# training

input:  $x$ , data:  $y$  noise:  $n$ .

$$d = y + n \Rightarrow n \sim \mathcal{N}(0, \sigma^2 I) \quad \text{즉, } \sigma \text{ 크기에 영향을 받는 것임.}$$

따라서,  $D_\theta$  direct  $x \Rightarrow D_\theta$ 를 제한 +  $F_\theta$  학습.

$$D_\theta(x; \sigma) = d - \sigma F_\theta(\cdot)$$

문제) large  $\sigma$  인 경우,  $\rightarrow$  Control,  $\rightarrow$   $F_\theta$ 가 너무 어려움.

본 논문에서 제안하는 것.

$$D_\theta(x; \sigma) = \underbrace{C_{skip}(\sigma)}_{\text{skip Connection}} x + \underbrace{C_{out}(\sigma)}_{\text{하위 모듈}} \underbrace{F_\theta}_{\text{scaler}}(\underbrace{C_{in}(\sigma)}_{\text{noise level}} x; C_{noise}(\sigma))$$

$$E_{\sigma, y, n} [\lambda(\sigma) \| D(y+n; \sigma) - y \|^2]$$

$$F_\theta \Rightarrow E_{\lambda, n} \left[ \underbrace{\lambda(\sigma) C_{out}(\sigma)^2}_{\text{effective weight}} \left\| \underbrace{F_\theta(C_{in}(\sigma) \cdot (y+n); C_{noise}(\sigma))}_{\text{network output}} - \underbrace{\frac{1}{C_{out}(\sigma)} (y - C_{skip}(\sigma) \cdot (y+n))}_{\text{training target}} \right\|^2 \right]$$

$$- \lambda(\sigma): \frac{1}{C_{out}(\sigma)^2} \Rightarrow \text{effective weight} = 1$$

$$- \text{noise level } \sigma \sim p_{\text{train}}(\sigma)$$

$\rightarrow$  너무 작은 경우  $\rightarrow$  너무 큰 경우  $\rightarrow$  normal distribution.



## Augmentation

→ overfitting 방지. 학습률 낮음.

본 레포트는 GAN의 Augmentation pipeline에