

latent space model SGM.

1. 새로운 score-matching objective.

2. 새로운 parameterization

3. 새로운 technic

1. Introduction

likelihood의 경우 data 분포를 잘 학습하면서도, high-fidelity인 data를 생성하는 것이 고민이 있다.

SGM은 data의 score function이 고려되었을 NCSN++와 SDE가 maximum likelihood로 바꾸는지를 살펴보자.

SGM은 cost가 많이 드는데, 간단한 noise 분포이며, 복잡한 분포로, SDE 및 probability flow ODE를 풀기 때문이다.

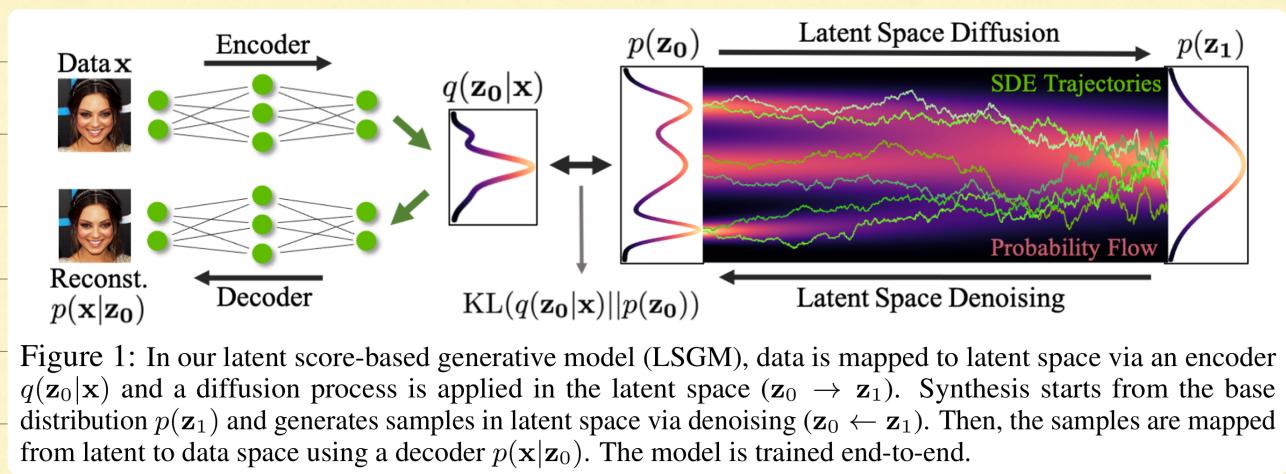
마지막 방정식은 복잡하고, 정확하게 풀려면,步를 작게 나누어야 한다는 때문에 소신개의 neural evaluate가 필요하다.

generation complexity: data distribution과 forward SDE의 차이로 결정되며, → 단계가 적으면 성능↑

SDE-based model은 continuous와 binary or categorical graph-structure data와 따라 적용성이 좋다.

LSGM은 VAE(?)와 latent score-based를 결합한다.

embedding space 생성 후 SGM은 embedding 분포학습



* Synthesis Speed.

VAE는 우선 gaussian의 prior를 pre-train 시킨다.

↳ prior가 diffusion의 noise인 비슷하도록 \Rightarrow diff. 속도.

* Expressivity

SGM은 neural ODE로 생각할 수 있다.

SGM은 latent space framework로 상상할 수 있다.

* Tailored Encoders and Decoders.

latent space의 SGM을 사용하여 encoder의 enc, decoder representation을 증가시킬 수 있다.

LSGM은 ELBO로 훈련 가능

일반 score-matching의 경우 DSM의 target distribution이 원래의 확률도가 되어야 물리적 확률과

기여 (contribution)

1. VAE의 LSGM을 동시에 확장시키는 DSM objective를 정의.

2. normal distribution과 learnable SGM의 경우 latent space function의 새로운 parameterization을

도입함, SGM의 latent의 분포와 normal prior 사이의 불일치만 modeling 할 수 있게 한다.

3. 새로운 SDE特殊的 importance sampling scheme을 도출함으로서, variance reduction tech로

안정적으로 훈련할 수 있도록 한다.

* Background.

NCSN $\leftrightarrow \frac{\nabla}{\lambda} \rightarrow \text{DSM}$.

\mathbf{z}_0 is start data, \mathbf{z}_t is time t'nes perturbed data, $t \in [0, 1]$ is continuous time step.

forward diffusion process $\{\mathbf{z}_t\}_{t=0}^{t=1}$ is

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\lambda(t) \mathbb{E}_{q(\mathbf{z}_0)} \mathbb{E}_{q(\mathbf{z}_t|\mathbf{z}_0)} [\|\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log p_{\theta}(\mathbf{z}_t)\|_2^2] \right] \quad (2) \rightarrow \text{SM}$$

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\lambda(t) \mathbb{E}_{q(\mathbf{z}_0)} \mathbb{E}_{q(\mathbf{z}_t|\mathbf{z}_0)} [\|\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t|\mathbf{z}_0) - \nabla_{\mathbf{z}_t} \log p_{\theta}(\mathbf{z}_t)\|_2^2] + C \right] \quad (3) \rightarrow \text{DSM}$$

$$\text{KL}(q(\mathbf{z}_0) \| p_{\theta}(\mathbf{z}_0)) \leq \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\frac{g(t)^2}{2} \mathbb{E}_{q(\mathbf{z}_0)} \mathbb{E}_{q(\mathbf{z}_t|\mathbf{z}_0)} [\|\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log p_{\theta}(\mathbf{z}_t)\|_2^2] \right] \quad (4)$$

$\hookrightarrow \text{DSM} \stackrel{?}{=} [12] \text{의 } \text{DDIM} \text{ 흐름}.$

* Score-based Generative modeling in Latent space.

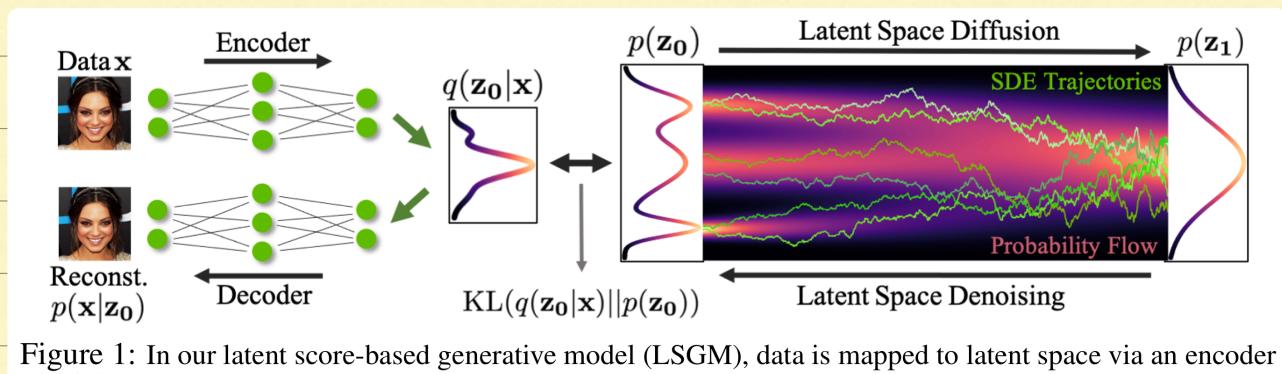


Figure 1: In our latent score-based generative model (LSGM), data is mapped to latent space via an encoder $q(\mathbf{z}_0|\mathbf{x})$ and a diffusion process is applied in the latent space ($\mathbf{z}_0 \rightarrow \mathbf{z}_1$). Synthesis starts from the base distribution $p(\mathbf{z}_1)$ and generates samples in latent space via denoising ($\mathbf{z}_0 \leftarrow \mathbf{z}_1$). Then, the samples are mapped from latent to data space using a decoder $p(\mathbf{x}|\mathbf{z}_0)$. The model is trained end-to-end.

Encoder: ϕ

Diffusion: $\Theta \in [0,1]$, $\mathbf{z}_t = \mathcal{N}(\mathbf{z}_0; 0, I)$

Decoder: ψ

Generate: Sampling from $\mathbf{z}_0 (\nabla_{\mathbf{z}_t} \log p_{\theta}(\mathbf{z}_t) \stackrel{?}{=} \nabla_{\mathbf{z}_t} \psi)$ then decoder ψ is used.

$$P(\mathbf{z}_0, \mathbf{x}) = P_{\theta}(\mathbf{z}_0) P_{\psi}(\mathbf{x}|\mathbf{z}_0)$$

훈련은 NLL의 ELBO로.

$$\mathcal{L}(\mathbf{x}, \phi, \theta, \psi) = \mathbb{E}_{q_\phi(\mathbf{z}_0|\mathbf{x})} [-\log p_\psi(\mathbf{x}|\mathbf{z}_0)] + \text{KL}(q_\phi(\mathbf{z}_0|\mathbf{x})||p_\theta(\mathbf{z}_0)) \quad (5)$$

$$= \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}_0|\mathbf{x})} [-\log p_\psi(\mathbf{x}|\mathbf{z}_0)]}_{\text{reconstruction term}} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}_0|\mathbf{x})} [\log q_\phi(\mathbf{z}_0|\mathbf{x})]}_{\text{negative encoder entropy}} + \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}_0|\mathbf{x})} [-\log p_\theta(\mathbf{z}_0)]}_{\text{cross entropy}} \quad (6)$$

어려운 점은 SGM의 CE 부분이다.

- The Cross-Entropy Term.

식(5)와 식(6)를 바로 KL로 대체시키지 않는 이유는 바로 계산이 불가능한데 때문이다.

$\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) \geq$ normalizing flow 같은 혼란적인 non-normal distribution의 경우 분석이 불가능하다.

DSM 또한 $\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) \geq C$ 의 형태로 나타나기 때문에 쉽게 놓거나 떼 수 없다.

이런 점을 극복하기 위해 eq.5를 decompose.

Theorem 1.

주어진 두 분포 $q(\mathbf{z}_0|\mathbf{x}), p(\mathbf{z}) \in \mathbb{R}^D$ 는 시간 티머의 $q(\mathbf{z}_t|\mathbf{x}), p(\mathbf{z}_t)$ 의 marginal distribution

$\log q(\mathbf{z}_t|\mathbf{x})$ 와 $\log p(\mathbf{z}_t)$ 를 smooth condition을 가정하면.

$$CE(q(\mathbf{z}_0|\mathbf{x})||p(\mathbf{z}_0)) = \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\frac{g(t)^2}{2} \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0|\mathbf{x})} \left[\|\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t|\mathbf{z}_0) - \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)\|_2^2 \right] \right] + \frac{D}{2} \log (2\pi e \sigma_0^2),$$

$$\hookrightarrow q(\mathbf{z}_t, \mathbf{z}_0|\mathbf{x}) = q(\mathbf{z}_t|\mathbf{z}_0)q(\mathbf{z}_0|\mathbf{x})$$

Normal transition kernel $q(\mathbf{z}_t|\mathbf{z}_0) = N(\mathbf{z}_t; \mu_t(\mathbf{z}_0), \Sigma_t^*)$

\downarrow
obtain from $f(t), g(t)$

eq.4와 다르게, $\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t)$ 이 유행 x

\hookrightarrow 증명 : Appendix A.

\hookrightarrow DSM 사용 가능 - $p(\mathbf{z}_0)$ optimally $q(\mathbf{z}_0|\mathbf{x})$ encoding으로 사용 가능

\hookrightarrow 복잡한 분포에서도 가능

Theorem (DLM) $\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$ 은 $p(\mathbf{z}_t)$ 분포를 diffusion score로 한다.

- Mixing Normal and Neural Score Functions

VAE나, 최근 사용되는 prior는 $N(\mathbf{z}_0; \mathbf{0}, \mathbf{I})$, prior를 $N(\mathbf{z}_t; \mathbf{0}, \mathbf{I})$ 로 가정한다.

single-D latent space를 가정하면, time t mixed mixture는

$$p(\mathbf{z}_t) \propto N(\mathbf{z}_t; \mathbf{0}, \mathbf{I})^{1-\alpha} \underbrace{p'_\theta(\mathbf{z}_t)^\alpha}_{\rightarrow \text{trainable SGIMU prior, } \alpha: \text{learnable scalar.}}$$

이점:

1) VAE pretrain 시킬 때, $\alpha=0$ 으로 설정하면 됨

pretrain의 prior는 $N(\mathbf{z}_0; \mathbf{0}, \mathbf{I})$ 로 만든다. SGIMU의 simple distribution을 학습하기 됨.

2) mixture의 score function은

$$\nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) = -((1-\alpha)\mathbf{z}_t + \alpha \nabla_{\mathbf{z}_t} \log p'_\theta(\mathbf{z}_t))$$

\hookrightarrow α 가 작을 때 linear한 미분하기 때문에 reverse SDE가 별리 좋음.

$\mathbf{z}_t \sim p(\mathbf{z}_t | \mathbf{z}_0)$ ($\mathbf{z}_t = \mu_t(\mathbf{z}_0) + \sigma_t \epsilon$ 일 때),

$$\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) = -\epsilon / \sigma_t \quad (\cong \text{DDPM}) \Rightarrow ?$$

$$\text{score function } \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) := -\epsilon_\theta(\mathbf{z}_t, t) / \sigma_t$$

$$\hookrightarrow \epsilon_\theta(\mathbf{z}_t, t) := \sigma_t(1-\alpha)\mathbf{z}_t + \alpha \circ \epsilon'_\theta(\mathbf{z}_t, t)$$

\downarrow CE 계산법

$$\text{CE}(q_\phi(\mathbf{z}_0 | \mathbf{x}) || p_\theta(\mathbf{z}_0)) = \mathbb{E}_{t \sim \mathcal{U}[0, 1]} \left[\frac{w(t)}{2} \mathbb{E}_{q_\phi(\mathbf{z}_t, \mathbf{z}_0 | \mathbf{x}), \epsilon} \left[\|\epsilon - \epsilon_\theta(\mathbf{z}_t, t)\|_2^2 \right] \right] + \frac{D}{2} \log \left(2\pi e \sigma_0^2 \right),$$

$$\hookrightarrow w(t) = g(t) / \sigma_t^2$$

- Training with different weight Mechanism

$w(t)$ 의 weight mechanism은 3가지 정도 있다.

Table 1: Weighting mechanisms

Mechanism	Weights	
Weighted	$w_{ll}(t) = g(t)^2 / \sigma_t^2$	maximum likelihood
Unweighted	$w_{un}(t) = 1$	DDPM
Reweighted	$w_{re}(t) = g(t)^2$	reweight

encoder $q(z_0|x)$ 은 $p(z_0|x)$ 의 1to1 대응으로 학습된다.

모든 Unweighted 비슷한 결과.

variance reduction이 더 간단한데, 두 가지 방법.

Training Objective

$$\min_{\phi, \psi} \mathbb{E}_{q_\phi(z_0|x)} [-\log p_\psi(x|z_0)] + \mathbb{E}_{q_\phi(z_0|x)} [\log q_\phi(z_0|x)] + \mathbb{E}_{t, \epsilon, q(z_t|z_0), q_\phi(z_0|x)} \left[\frac{w_{ll}(t)}{2} \|\epsilon - \epsilon_\theta(z_t, t)\|_2^2 \right] \quad (8)$$

$$\min_{\theta} \mathbb{E}_{t, \epsilon, q(z_t|z_0), q_\phi(z_0|x)} \left[\frac{w_{ll/un/re}(t)}{2} \|\epsilon - \epsilon_\theta(z_t, t)\|_2^2 \right] \quad \text{with } q(z_t|z_0) = \mathcal{N}(z_t; \mu_t(z_0), \sigma_t^2 \mathbf{I}), \quad (9)$$

- Variance Reduction.

eq 8.9는 high variance를 가진다.

넓은 줄이다.

Variance preserving : $dz = -\frac{1}{2} \beta(t) z dt + \sqrt{\beta(t)} dw$, $\beta(t) = \beta_0 + (\beta_1 - \beta_0)t$.

- Variance reduction for likelihood weighting

Appendix B) M. $q(z_0) = p(z_0) = N(z_0; 0, I)$ 일 때,

$$CE(q(z_0) \| p(z_0)) = \frac{D}{2} \mathbb{E}_{t \sim U[0, 1]} [d \log \sigma_t^2 / dt] + \text{const.}$$

1) Geometric VPSDE.

uniform t 의 case variance가 줄어들고 있음, $d \log \sigma_t^2 / dt \rightarrow t \sim [0, 1]$ 의 경우 uniform

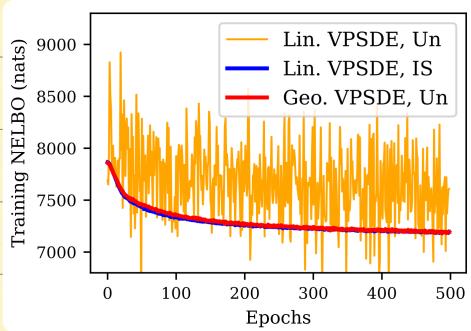
$$\beta(t) = \log(\sigma_{\max}^2 / \sigma_{\min}^2) \frac{\sigma_t^2}{(1 - \sigma_t^2)} \quad \sigma_t^2 = \sigma_{\min}^2 (\sigma_{\max}^2 / \sigma_{\min}^2)^t \text{ 일 때 만족.}$$

VEISPEL의 경우, t 의 case reduction variance는 20%

2) Importance Sampling (IS)

$$r(t) \propto d \log \sigma_t^2 / dt \text{ 만족}$$

$$t = \text{var}^{-1}((\sigma_t^2)^p (\sigma_0^2)^{1-p})$$



- Variance reduction for Unweight and reweight.

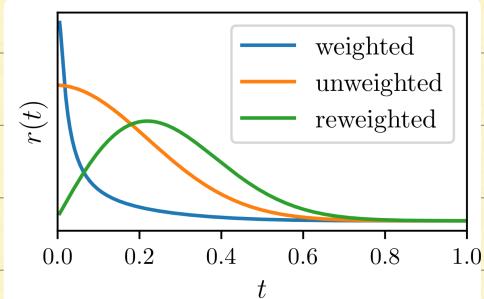


Figure 3: IS distributions

App B, G

* Related Work.

* Experiments.

NVAE + NCSN++

black box ODE solver for sampling

Table 2: Generative performance on CIFAR-10.

	Method	NLL↓	FID↓
Ours	LSGM (FID)	≤ 3.43	2.10
	LSGM (NLL)	≤ 2.87	6.89
	LSGM (balanced)	≤ 2.95	2.17
	VAE Backbone	2.96	43.18
VAEs	VDVAE [21]	2.87	-
	NVAE [20]	2.91	23.49
	VAEBM [76]	-	12.19
	NCP-VAE [56]	-	24.08
	BIVA [48]	3.08	-
	DC-VAE [77]	-	17.90
Score	NCSN [3]	-	25.32
	Rec. Likelihood [40]	3.18	9.36
	DSM-ALS [39]	3.65	-
	DDPM [1]	3.75	3.17
	Improved DDPM [26]	2.94	11.47
	SDE (DDPM++) [2]	2.99	2.92
Flows	SDE (NCSN++) [2]	-	2.20
	VFlow [19]	2.98	-
	ANF [18]	3.05	-
Aut. Reg.	DistAug aug [78]	2.53	42.90
	Sp. Transformers [79]	2.80	-
	δ -VAE [80]	2.83	-
	PixelSNAIL [81]	2.85	-
	PixelCNN++ [82]	2.92	-
GANs	AutoGAN [83]	-	12.42
	StyleGAN2-ADA [84]	-	2.92

Table 4: Dyn. binarized OMNIGLOT results.

	Method	NELBO↓	NLL↓
Ours	LSGM	87.79	≤ 87.79
VAEs	NVAE [20]	93.92	90.75
	BIVA [48]	93.54	91.34
	DVAE++ [51]	-	92.38
	Ladder VAE [90]	-	102.11
Aut. Reg.	VLVAE [47]	-	89.83
	VampPrior [59]	-	89.76
	PixelVAE++ [91]	-	88.29

Table 3: Generative results on CelebA-HQ-256.

	Method	NLL↓	FID↓
Ours	LSGM	≤ 0.70	7.22
	VAE Backbone	0.70	30.87
VAEs	NVAE [20]	0.70	29.76
	VAEBM [76]	-	20.38
	NCP-VAE [56]	-	24.79
	DC-VAE [77]	-	15.80
Score	SDE [2]	-	7.23
Flows	GLOW [85]	1.03	68.93
Aut. Reg.	SPN [86]	0.61	-
GANs	Adv. LAE [87]	-	19.21
	VQ-GAN [64]	-	10.70
	PGGAN [88]	-	8.03

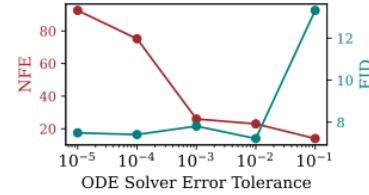


Figure 4: FID and number of function evaluations (NFEs) for different ODE solver error tolerances on CelebA-HQ-256. LSGM takes 4.15 sec. for sampling while the original SGM [2] takes 45 min. with PC and 3.9 min. with ODE-based sampling.

Table 5: Dynamically binarized MNIST results.

	Method	NELBO↓	NLL↓
Ours	LSGM	78.47	≤ 78.47
VAEs	NVAE [20]	79.56	78.01
	BIVA [48]	80.06	78.41
	IAF-VAE [24]	80.80	79.10
	DVAE++ [51]	-	78.49
Aut. Reg.	PixelVAE++ [91]	-	78.00
	VampPrior [59]	-	78.45
	MAE [92]	-	77.98



(a) CIFAR-10

(b) CelebA-HQ-256

(c) OMNIGLOT

1 0 8 7 1 6 4 8
2 0 8 7 2 5 0 5
6 2 3 0 7 3 6 7
7 5 0 5 5 7 9 8

(d) MNIST

Figure 5: Generated samples for different datasets. For binary datasets, we visualize the decoder mean. LSGM successfully generates sharp, high-quality, and diverse samples (additional samples in appendix).

Table 6: Ablations on SDEs, objectives, weighting mechanisms, and variance reduction. Details in App. G.

SGM-obj.-weighting		w_{II}		w_{un}				w_{re}			
<i>t</i> -sampling (SGM-obj.)		$\mathcal{U}[0, 1]$	$r_{\text{II}}(t)$	$\mathcal{U}[0, 1]$		$r_{\text{un}}(t)$	$\mathcal{U}[0, 1]$		$r_{\text{re}}(t)$		
Geom.- VPSDE	$FID \downarrow$ $NELBO \downarrow$	rew.	rew.	rew.	$r_{\text{II}}(t)$	rew.	$r_{\text{II}}(t)$	rew.	$r_{\text{II}}(t)$	rew.	$r_{\text{II}}(t)$
Geom.- VPSDE	$FID \downarrow$ $NELBO \downarrow$	10.18 2.96	n/a n/a	NaN NaN	NaN NaN	n/a n/a	n/a n/a	22.21 3.04	NaN NaN	7.29 2.99	7.18 2.99
VPSDE	$FID \downarrow$ $NELBO \downarrow$	6.15 2.97	8.00 2.97	NaN NaN	NaN NaN	5.39 2.98	5.39 2.98	NaN NaN	4.99 2.99	15.12 3.03	6.19 2.99

* Conclusion.

latent diff diffusion model

↳ more representation.

↳ discrete image modeling, smooth SGM.

Contribution.

1) MSEL CE objective 조합

2) Mixing Prior & prior parameterization

3) Variance reduction

SGM for t-sampling