

Generative models for when inference gets complicated.

Gabriel V. Cardoso and Mike Pereira (Ecole des Mines - PSL)



GEOLEARNING
CHAIR Data Science for the Environment



INRAE
la science pour la vie, l'ambiance, le territoire



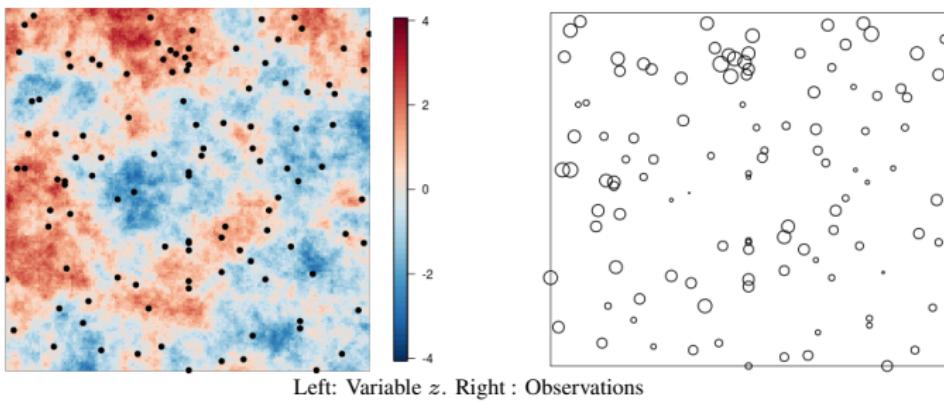
SCOR
FONDATION POUR LA SCIENCE

April 14, 2025

The problem

Motivation

- Data: Observations $\{x(s_i)\}_{i=1}^n$ of a variable x at n locations s_1, \dots, s_n of a domain $\mathcal{D} \subset \mathbb{R}^d$



→ Goal: Predict x at a new location $s_0 \in \mathcal{D}$

Motivation

Geostatistical paradigm: over the spatial domain \mathcal{D}

Gaussian Random Field

$$X : \{X(s) : s \in \mathcal{D}\}$$

Realization
 $\xrightarrow{\hspace{1cm}}$

High correlation

Observed variable

$$x : \{x(s) : s \in \mathcal{D}\}$$

High *similarity*

- Finite-dimensional distributions fully specified by mean function + covariance function \rightarrow Parametrized by θ
- Analytic formulas for predictive posterior distribution

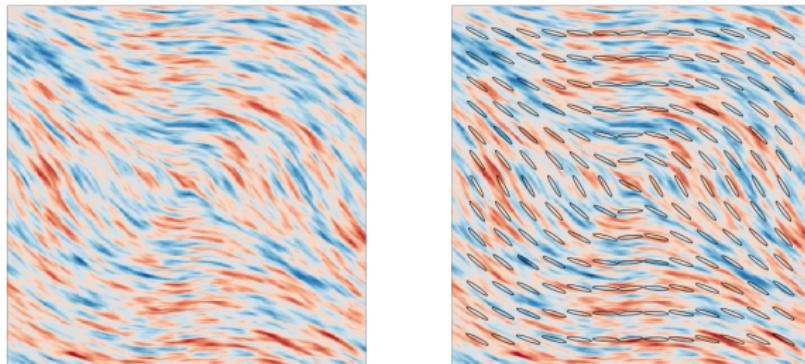
$$(X(s_0) | X(s_1), \dots, X(s_n), \theta)$$

$\sim \mathcal{N}(\text{Kriging prediction at } s_0, \text{Kriging variance at } s_0)$

\rightarrow BUT Need to **estimate parameters** θ from data to compute it!

Motivation

- Parameter estimation relatively straightforward for isotropic models (Max-likelihood, MCMC,...)
- But what if we consider non-stationary Gaussian random fields (with locally varying parameters) ?



- Maps of parameters \Rightarrow Complex parameter space \Rightarrow Max-Likelihood/MCMC harder to converge
- INLA approach (i.e. Laplace approximation) implemented for less complex models (varying range, variance)

Motivation

- Idea : Use generative models to perform Bayesian predictions
- Our setting :
 - Set prior distribution on parameters θ
 - Use generative models to learn the corresponding prior distribution of GRFs
 - Use conditioning of generative model to sample from the posterior predictive distribution $(X(s_0)|X(s_1), \dots, X(s_n))$ (\rightarrow Marginalized over θ !)

Fast introduction to score-based generative models

The original miracle!



Figure: $X_0 \sim \pi$.



Figure: $X_\sigma = X_0 + \sigma Z$.¹ 1

¹ $Z \sim \mathcal{N}(0, I)$, $X_0 \perp Z$.

The original miracle!



Figure: $X_0 \sim \pi$.



Figure: $X_\sigma = X_0 + \sigma Z$.¹

- D_{θ_0} untrained CNN,
- $L_\sigma(\theta) = \|D_\theta(X_\sigma) - X_\sigma\|^2$, $L_0(\theta) = \|D_\theta(X_\sigma) - X_0\|^2$,
- $\theta_t = \theta_{t-1} - \nabla L_\sigma(\theta_{t-1})$

¹ $Z \sim \mathcal{N}(0, I)$, $X_0 \perp Z$.

The original miracle!



Figure: $X_0 \sim \pi$.



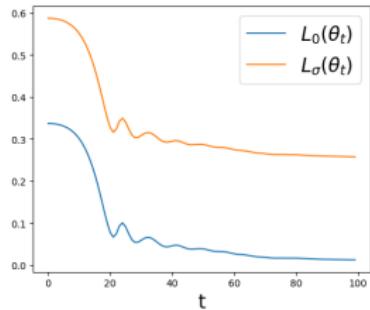
Figure: $X_\sigma = X_0 + \sigma Z$.¹ 1

- D_{θ_0} untrained CNN,
- $L_\sigma(\theta) = \|D_\theta(X_\sigma) - X_\sigma\|^2$, $L_0(\theta) = \|D_\theta(X_\sigma) - X_0\|^2$,
- $\theta_t = \theta_{t-1} - \nabla L_\sigma(\theta_{t-1})$

What happens to $\{L_\sigma(\theta_t)\}_{t=1}^N$ and $\{L_0(\theta_t)\}_{t=1}^N$?

¹ $Z \sim \mathcal{N}(0, I)$, $X_0 \perp Z$.

The original miracle ².



²This experiment is inspired by Ulyanov et al., “Deep Image Prior”, 2018

Noising and Denoising



Figure: $X_0 \sim \pi$.

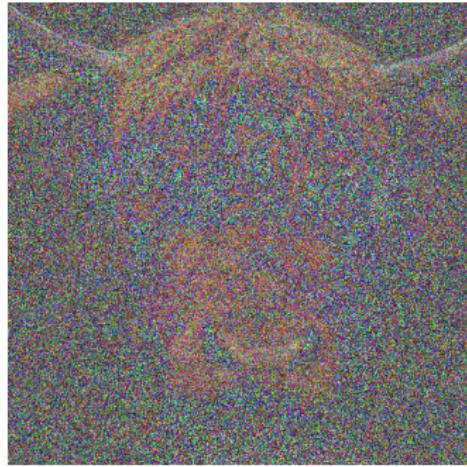


Figure: $X_t = X_0 + \sigma_t Z_t$.³¹

Suppose we have an increasing noise schedule

$$\sigma : [0, T] \rightarrow \sigma_t \in \mathbb{R}_+$$

Can we reverse the procedure?

¹ $Z_t \perp X_0, Z_t \sim \mathcal{N}(0, I)$.

Noising and Denoising

Suppose we have an **increasing noise schedule**

$$\sigma : [0, T] \rightarrow \sigma_t \in \mathbb{R}_+ .$$

1 Noising kernel:

$$p_{t|0}(x_t|x_0) = \mathcal{N}(x_t; x_0, \sigma_t^2 \mathbf{I}) ,$$

2 Marginal:

$$p_t(x_t) = \int p_{t|0}(x_t|x_0) \pi(\mathrm{d}x_0) .$$

Our goal is to go backwards:

$$p_T \rightarrow p_{T-\Delta t} \rightarrow \cdots \rightarrow p_0 .$$



Figure: $X_t = X_0 + \sigma_t Z_t$.
41

Generative models by score matching⁵

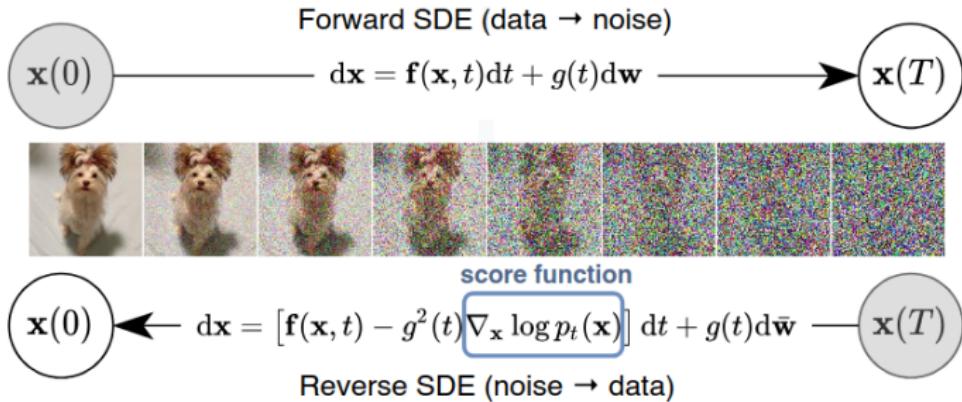


Figure: Illustration of noising and denoising with SDEs.

⁴From Song et al., “Score-Based Generative Modeling through Stochastic Differential Equations”, 2021.

⁵Song et al., “Score-Based Generative Modeling through Stochastic Differential Equations”, 2021.

Generative models by score matching⁵

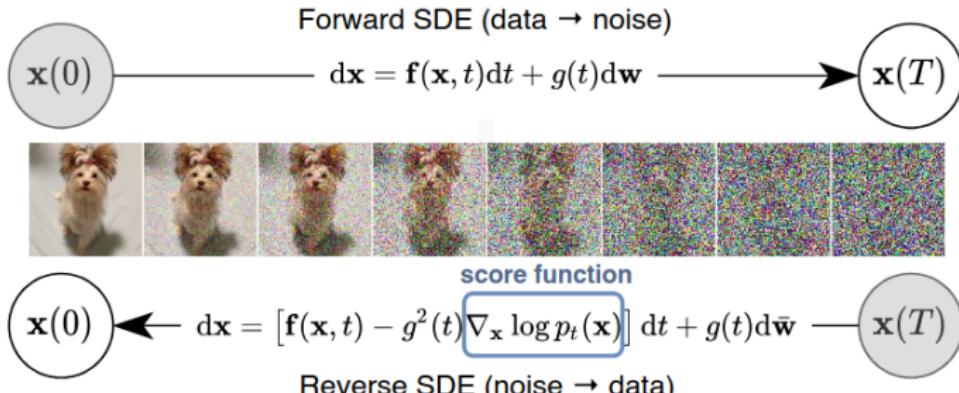


Figure: Illustration of noising and denoising with SDEs.

Our p_t is the time t marginal of the SDE

$$dx = dw .$$

It admits the backward decomposition

$$dx = -\nabla \log p_t(x)dt + d\bar{w} .$$

⁴From Song et al., “Score-Based Generative Modeling through Stochastic Differential Equations”, 2021.

⁵Song et al., “Score-Based Generative Modeling through Stochastic Differential Equations”, 2021.

Theoretical guarantees

We can control

- $\text{KL}(\mathbf{q}_{0:T}(x_{0:T}) \parallel \mathbf{p}_{0:T}(x_{0:T}))$ ⁶
- $\text{W}_2(\mathbf{p}_{0:T}(x_{0:T}), \mathbf{q}_{0:T}(x_{0:T})).$ ⁷

⁶Conforti et al., *KL Convergence Guarantees for Score Diffusion Models under Minimal Data Assumptions*, 2024.

⁷Gao et al., *Wasserstein Convergence Guarantees for a General Class of Score-Based Generative Models*, 2025.

Denoising: Optimal denoiser and the score

Consider the mean squared error (MSE)

$$\text{MSE} \left(f, \sigma_t^2 \right) := \mathbb{E}_{\mathbf{q}_{t,0}} [\|f(X_t) - X_0\|^2].$$

⁸Vincent, “A Connection Between Score Matching and Denoising Autoencoders”, 2011.

Denoising: Optimal denoiser and the score

Consider the mean squared error (MSE)

$$\text{MSE} \left(f, \sigma_t^2 \right) := \mathbb{E}_{\mathbf{q}_{t,0}} [\|f(X_t) - X_0\|^2].$$

The **optimal denoiser** w.r.t the MSE is,

$$x_t \rightarrow \mathbb{E} [X_0 | X_t = x_t] \in \operatorname*{argmin}_{f \in L^2(\mathbf{q}_t)} \text{MSE} \left(f, \sigma_t^2 \right).$$

⁸Vincent, “A Connection Between Score Matching and Denoising Autoencoders”, 2011.

Denoising: Optimal denoiser and the score

Consider the mean squared error (MSE)

$$\text{MSE} \left(f, \sigma_t^2 \right) := \mathbb{E}_{\mathbf{q}_{t,0}} [\|f(X_t) - X_0\|^2].$$

The **optimal denoiser** w.r.t the MSE is,

$$x_t \rightarrow \mathbb{E} [X_0 | X_t = x_t] \in \operatorname{argmin}_{f \in L^2(\mathbf{q}_t)} \text{MSE} \left(f, \sigma_t^2 \right).$$

It is possible to show that⁸

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \underbrace{\nabla \log \mathbf{q}_t(x_t)}_{\text{score}}.$$

⁸Vincent, “A Connection Between Score Matching and Denoising Autoencoders”, 2011.

Denoising: Optimal denoiser and the score

Consider the mean squared error (MSE)

$$\text{MSE} \left(f, \sigma_t^2 \right) := \mathbb{E}_{\mathbf{q}_{t,0}} [\|f(X_t) - X_0\|^2].$$

The **optimal denoiser** w.r.t the MSE is,

$$x_t \rightarrow \mathbb{E} [X_0 | X_t = x_t] \in \underset{f \in L^2(\mathbf{q}_t)}{\operatorname{argmin}} \text{MSE} \left(f, \sigma_t^2 \right).$$

It is possible to show that⁸

$$\mathbb{E} [X_0 | X_t = x_t] = x_t + \sigma_t^2 \underbrace{\nabla \log \mathbf{q}_t(x_t)}_{\text{score}}.$$

Learn a **denoiser** \iff Learn the **score**.

⁸Vincent, “A Connection Between Score Matching and Denoising Autoencoders”, 2011.

Backward sampling from Ho et al., “Denoising diffusion probabilistic models”, 2020

Algorithm Backward sampling

1: **procedure** SAMPLING

Input: N .

Output: X_0 .

2: $X_N \sim \mathcal{N}(0, \sigma_T^2 \mathbf{I})$.

3: **for** $i = N - 1, \dots, 0$ **do**

4: $t_i = (i/N)T$

5: $Z_i \sim \mathcal{N}(0, \mathbf{I})$.

6:
$$X_i = X_{i+1} + \frac{\sigma_{t_i}^2 - \sigma_{t_{i+1}}^2}{\sigma_{t_{i+1}}^2} (X_{i+1} - D_\theta(X_{i+1}, \sigma_{t_{i+1}})) + \sqrt{\frac{\sigma_{t_i}^2}{\sigma_{t_{i+1}}^2} (\sigma_{t_{i+1}}^2 - \sigma_{t_i}^2)} Z_i.$$

7: **end for**

8: **Return:** X_0

9: **end procedure**

DDPM: Illustration

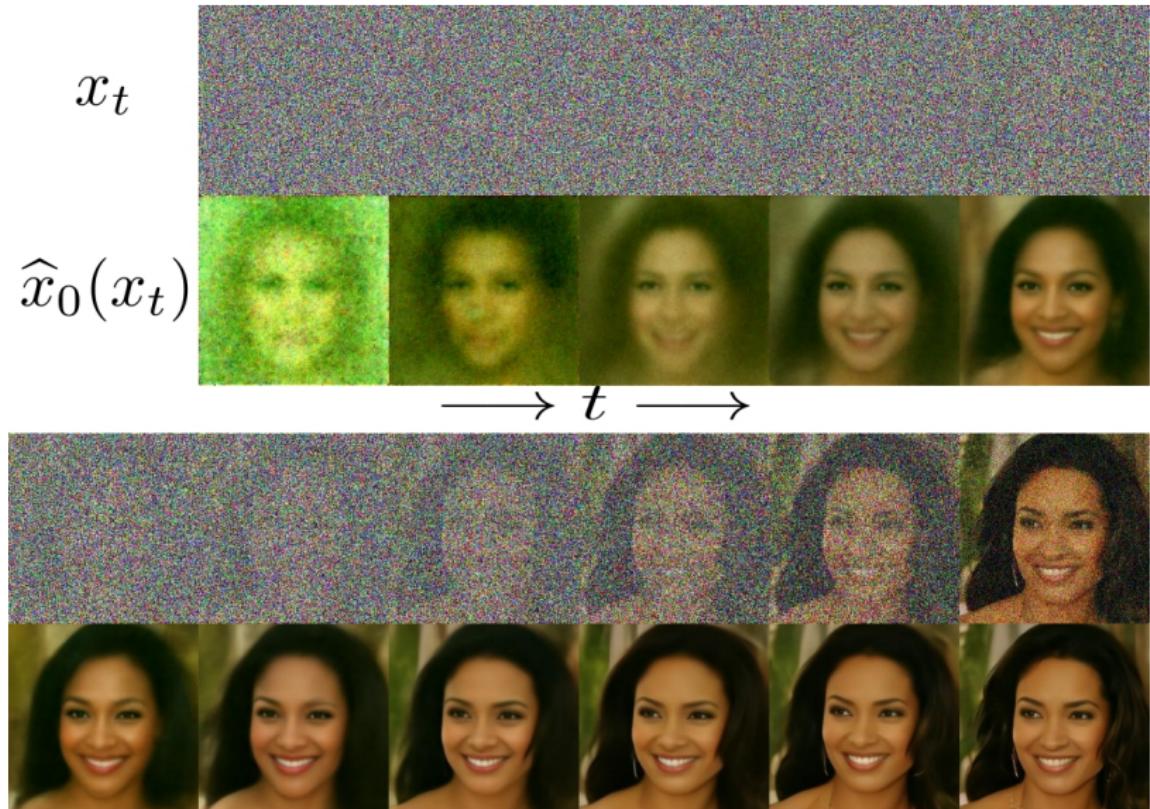


Figure: Illustration of DDPM on CelebA dataset.

Denoising with UNETs

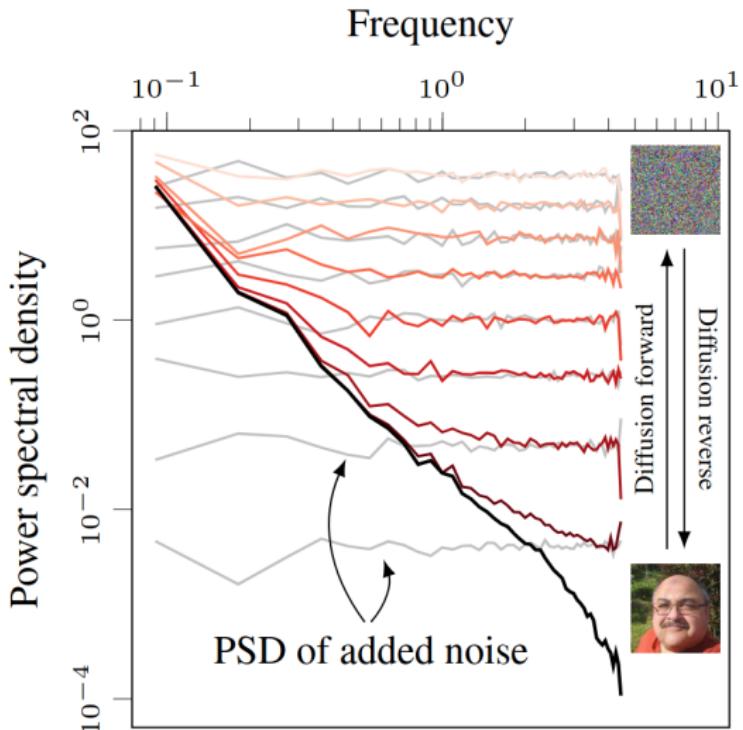


Figure: Illustration from Severi Rissanen et al. “Generative Modelling with Inverse Heat Dissipation”. In: *The Eleventh International Conference on Learning Representations*. Sept. 2022. (Visited on 11/13/2024).

Denoising with UNETs: Geometry-adapted harmonic representation

Linear denoising in an adaptative basis:

$$\hat{x}_0(x_t; \theta) = \operatorname{argmin}_{\lambda, \mathcal{B} \in \mathcal{D}} \sum_{e_k \in \mathcal{B}} \lambda_k(x_t) e_k(x_t)$$

Denoising with Unets in C^α images

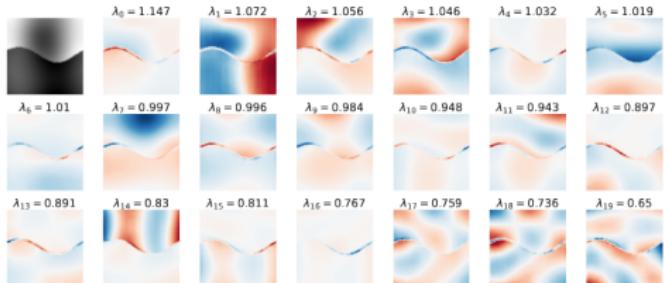
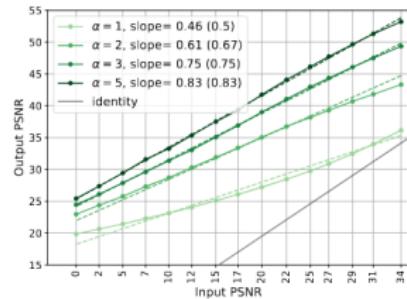


Figure: Illustration from Zahra Kadkhodaie et al. “Generalization in diffusion models arises from geometry-adaptive harmonic representation”. en. In: *The Twelfth International Conference on Learning Representations*. 2023.

$$\text{MSE}(f, \sigma) \propto \sigma^{2\frac{\alpha}{\alpha+1}}$$

Denoising with Unets: Celeb

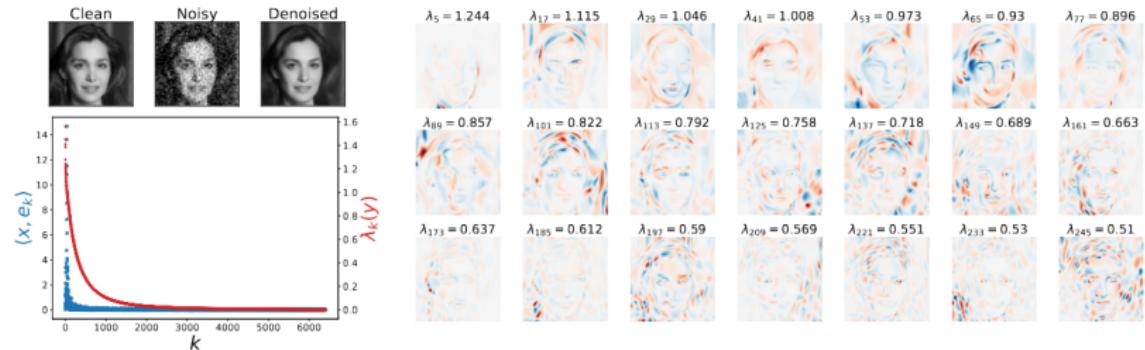


Figure: Illustration from Zahra Kadkhodaie et al. “Generalization in diffusion models arises from geometry-adaptive harmonic representation”. en. In: *The Twelfth International Conference on Learning Representations*. 2023.

Denoising with UNETs

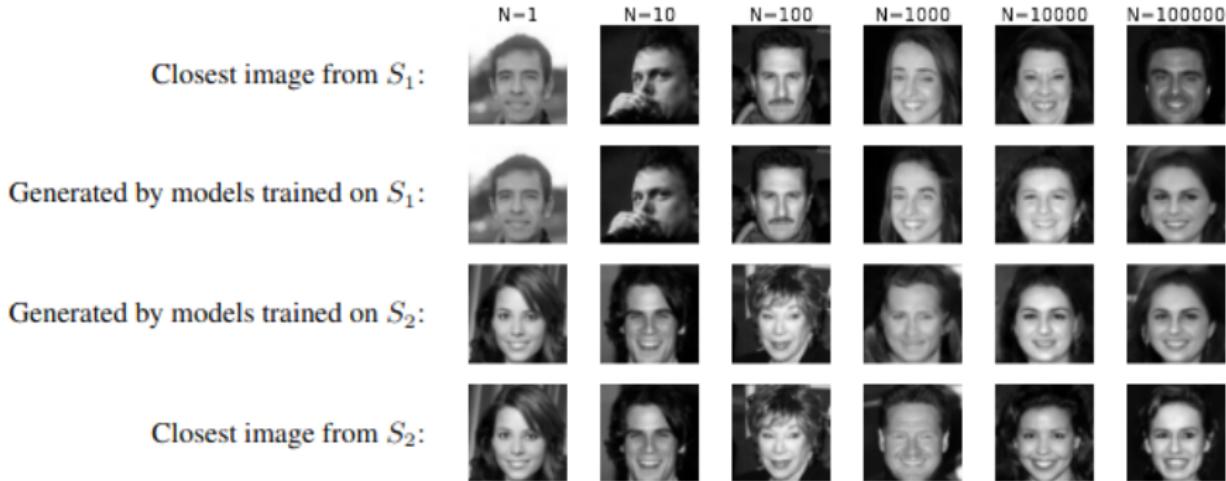


Figure: Illustration from Zahra Kadkhodaie et al. “Generalization in diffusion models arises from geometry-adaptive harmonic representation”. en. In: *The Twelfth International Conference on Learning Representations*. 2023.

Struggles and Results

Quizz

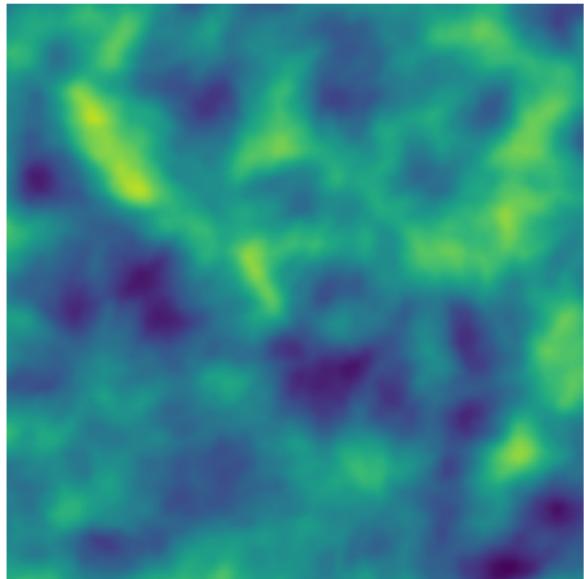


Figure: 0

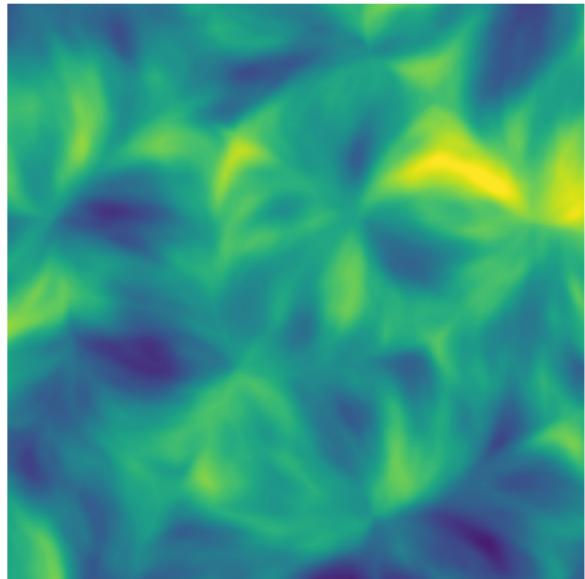


Figure: 1

Quizz

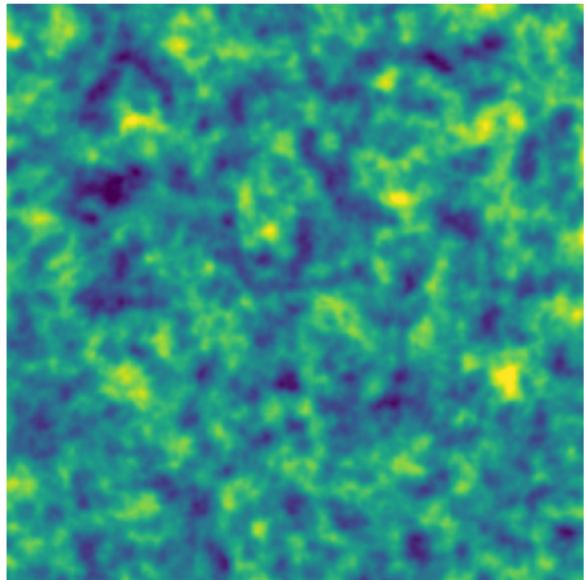


Figure: 0

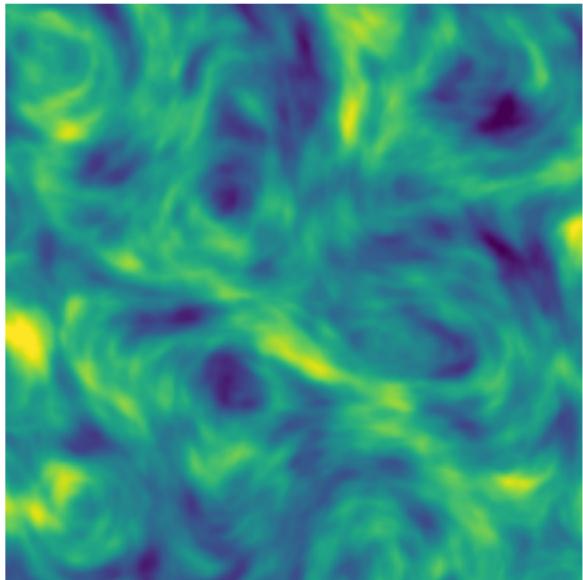


Figure: 1

Quizz

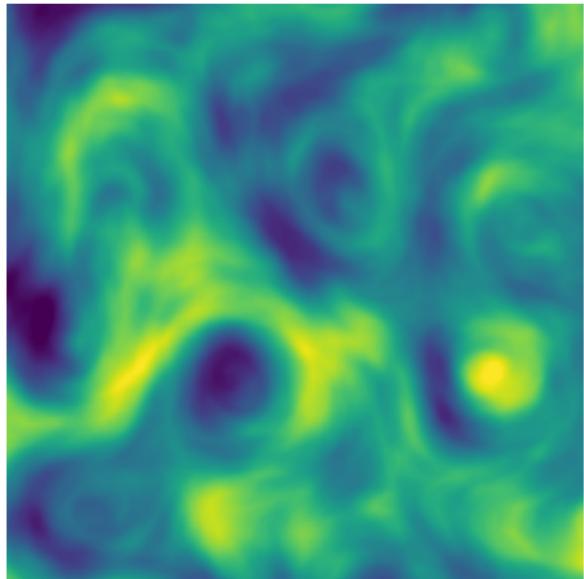


Figure: 0

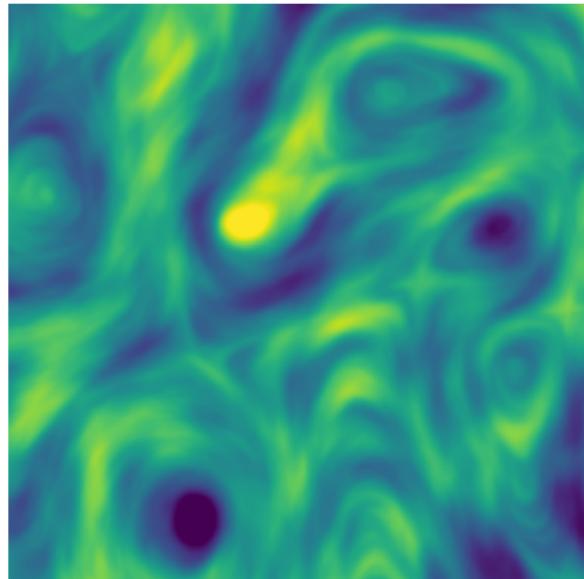


Figure: 1

Quizz

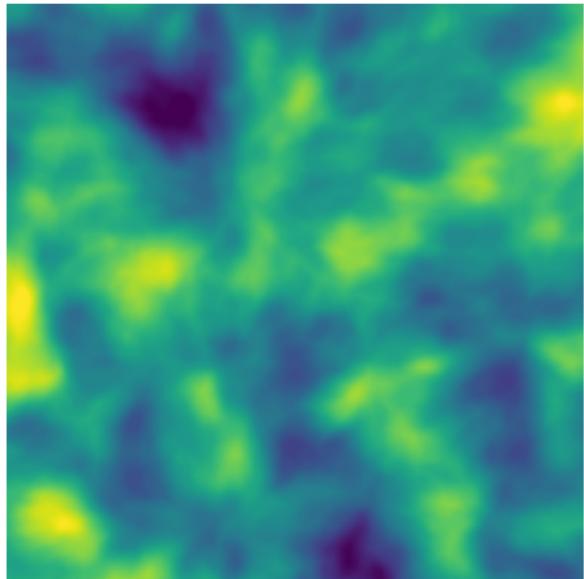


Figure: 0

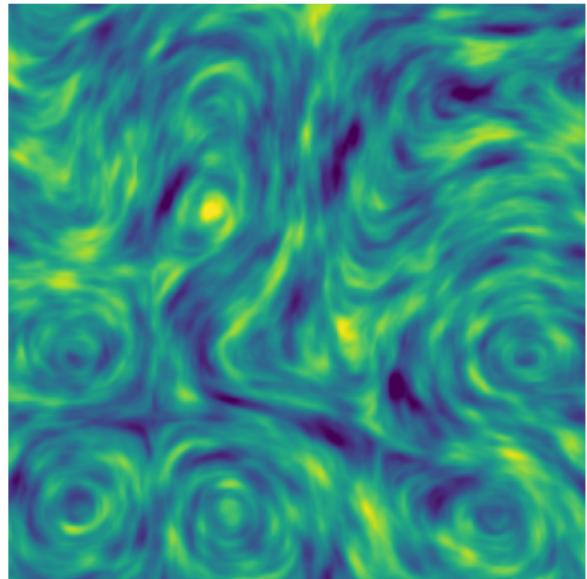


Figure: 1

Quizz

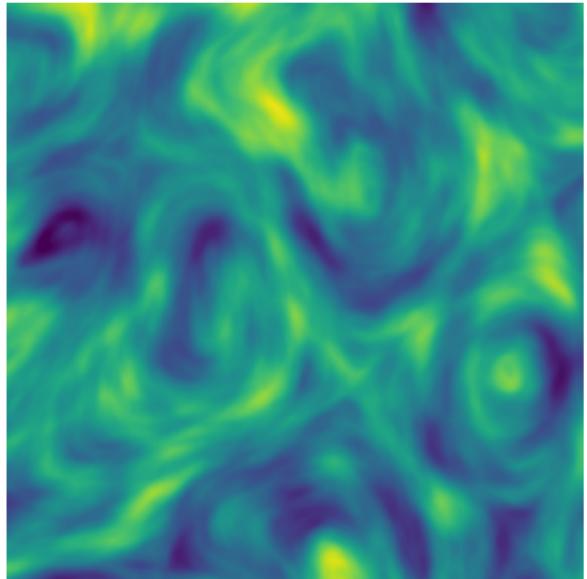


Figure: 0

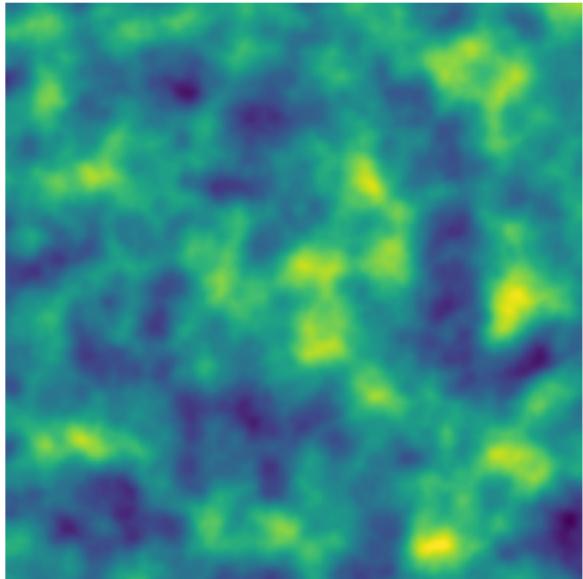


Figure: 1

Quizz

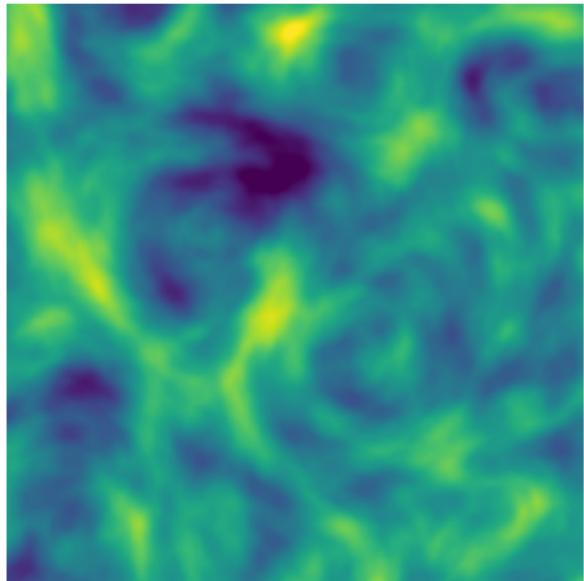


Figure: 0

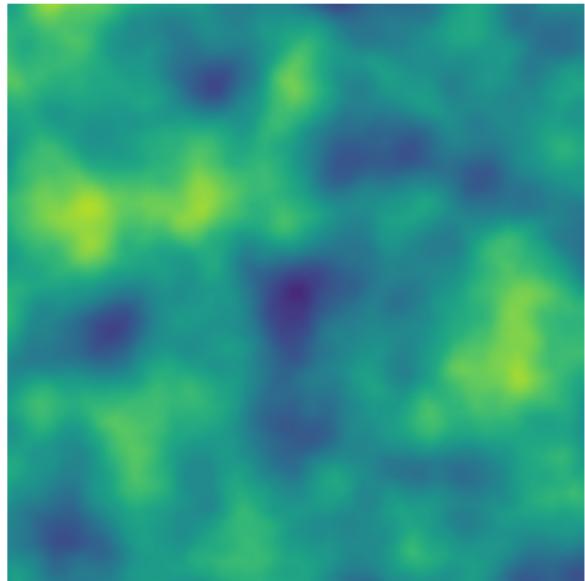


Figure: 1

Illustration

How to judge? Pt.1 Max-Sliced Wasserstein

Wasserstein distances

Let $\mu, \nu \in \mathcal{P}_p$. We define

$$W_p(\mu, \nu)^p := \inf_{\gamma \in C(\mu, \nu)} \int \|x - y\|^p \gamma(dx, dy), \quad (1)$$

⁹See Nietert et al., “Statistical, Robustness, and Computational Guarantees for Sliced Wasserstein Distances”, 2022 for guarantees.

How to judge? Pt.1 Max-Sliced Wasserstein

Wasserstein distances

Let $\mu, \nu \in \mathcal{P}_p$. We define

$$W_p(\mu, \nu)^p := \inf_{\gamma \in C(\mu, \nu)} \int \|x - y\|^p \gamma(dx, dy), \quad (1)$$

$$\overline{W}_p(\mu, \nu) = \max_{\theta \in \mathbb{S}^{d-1}} W_p(\mathfrak{p}^\theta \# \mu, \mathfrak{p}^\theta \# \nu). \quad (2)$$

⁹See Nietert et al., “Statistical, Robustness, and Computational Guarantees for Sliced Wasserstein Distances”, 2022 for guarantees.

How to judge? Pt.1 Max-Sliced Wasserstein

Wasserstein distances

Let $\mu, \nu \in \mathcal{P}_p$. We define

$$W_p(\mu, \nu)^p := \inf_{\gamma \in C(\mu, \nu)} \int \|x - y\|^p \gamma(dx, dy), \quad (1)$$

$$\overline{W}_p(\mu, \nu) = \max_{\theta \in \mathbb{S}^{d-1}} W_p(\mathfrak{p}^\theta \# \mu, \mathfrak{p}^\theta \# \nu). \quad (2)$$

Empirical Wasserstein ⁹

Suppose $X_{1:n} \sim \mu^{\otimes n}$, $Y_{1:n} \sim \nu^{\otimes n}$, and $\hat{\mu}_n(dx) = n^{-1} \sum_{i=1}^n \delta_{X_i}(dx)$ and $\hat{\nu}_n(dx) = n^{-1} \sum_{i=1}^n \delta_{Y_i}(dx)$. We define

$$\overline{W}_{MC,p}(\hat{\mu}_n, \hat{\nu}_n; K) := \max_{\theta \in \Theta} W_p(\mathfrak{p}^\theta \# \hat{\mu}_n, \mathfrak{p}^\theta \# \hat{\nu}_n), \quad (3)$$

where $\Theta \sim \text{Unif}(\mathbb{S}^{d-1})^{\otimes K}$.

⁹See Nietert et al., “Statistical, Robustness, and Computational Guarantees for Sliced Wasserstein Distances”, 2022 for guarantees.

How to judge? Pt.2 Classifier two-sample test

- Generate N samples from a given sampler,
- Create a classification dataset using the test set and the generated samples.
- Split in two
- Train on one, validate on other
- Check your favorite metric.

Results: Max-Sliced Wasserstein

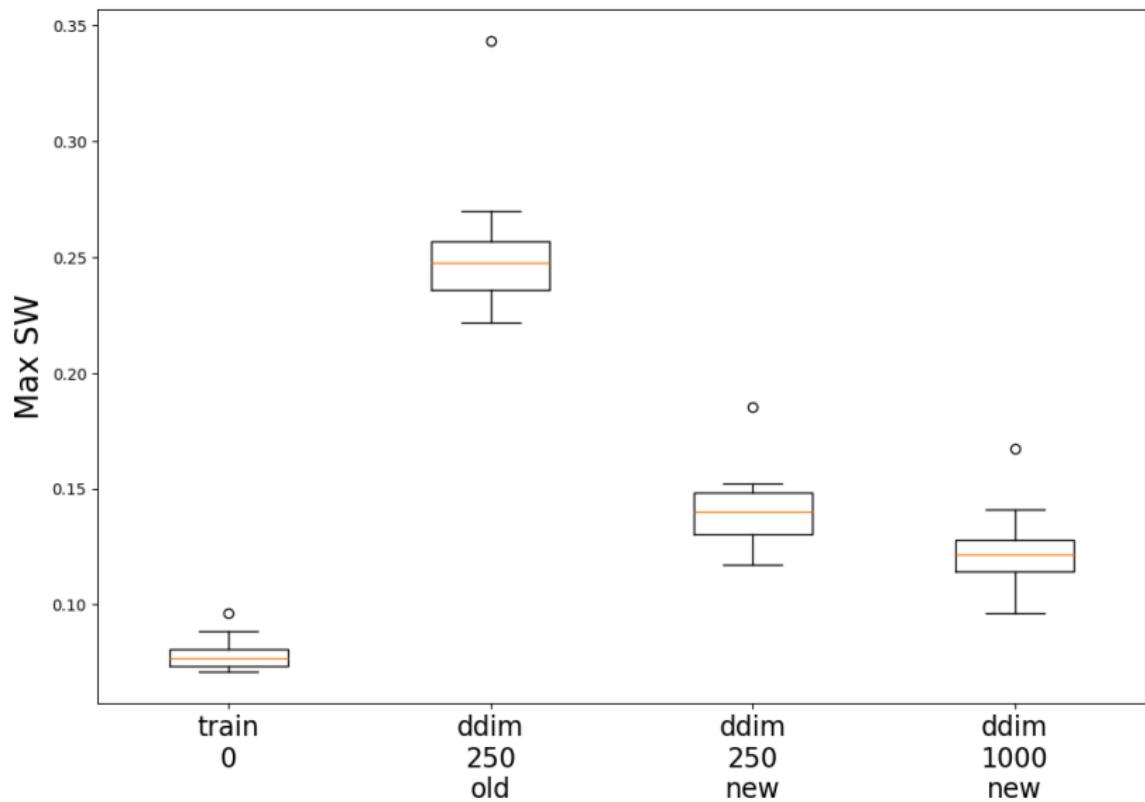


Figure: Max-Sliced Wasserstein

Results: Classifier (Resnet18)

sampler	AUC	Network
DDIM 250 old	0.84	Resnet18
DDIM 250 new	0.81	Resnet18
DDIM 1000 new	0.61	Resnet18
DDIM 1000 new	0.60	Resnet50

Table: ROC AUC.

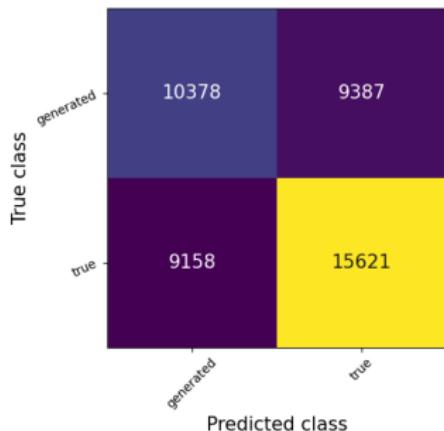
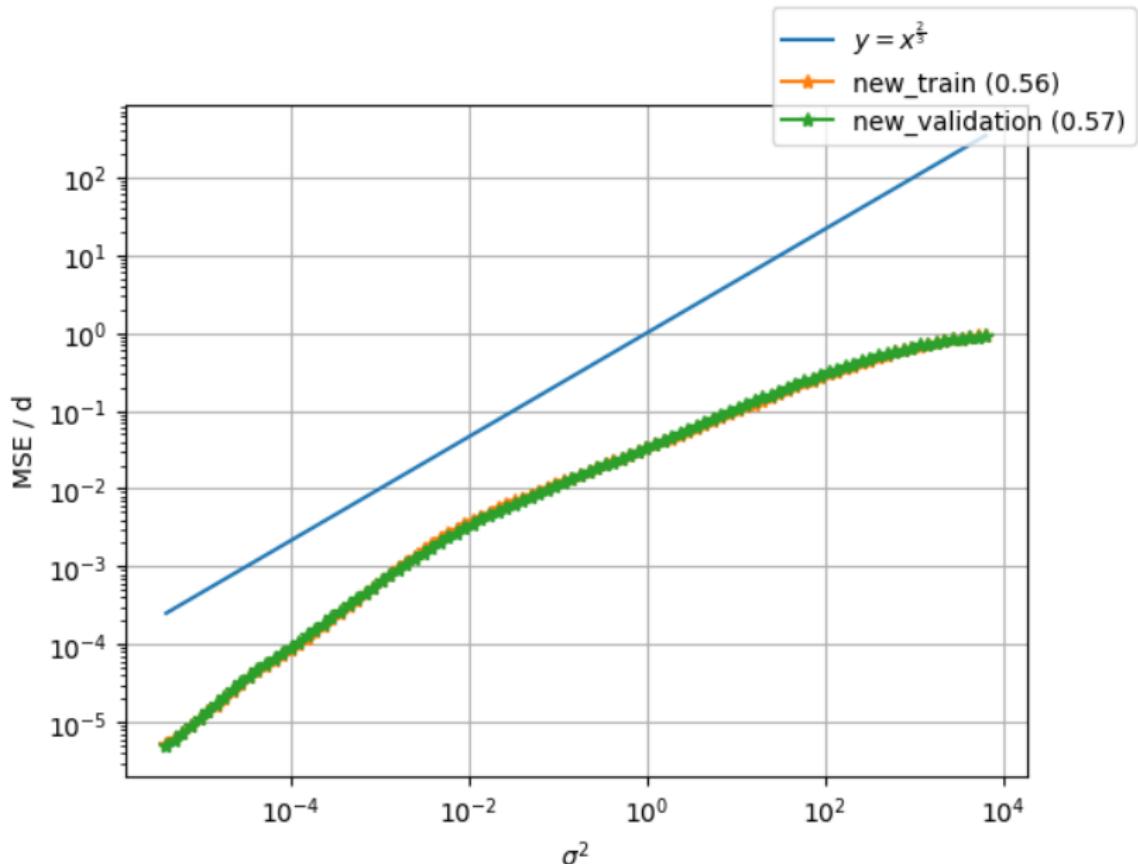


Figure: Confusion matrix for DDIM 1000 with ResNet18.

MSE vs σ^2



A glimpse of posterior sampling

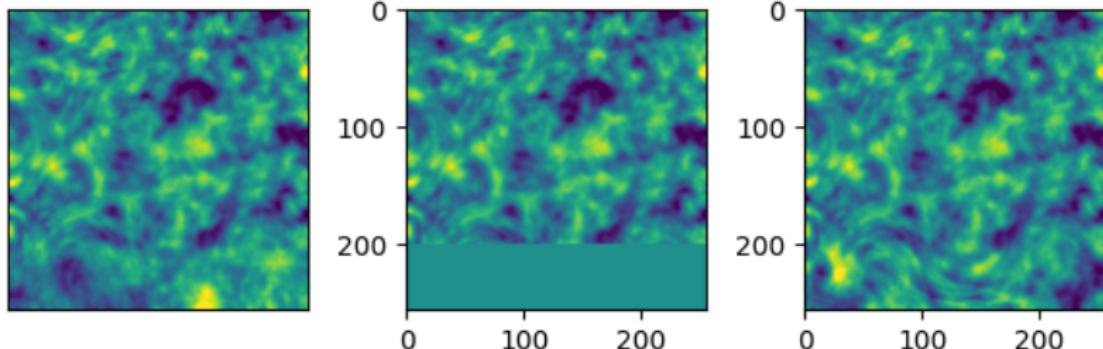


Figure: Posterior sampling result with old model.

Perspectives and challenges

- Should we treat different regularities?
- How to validate posterior sampling?