

Bayesian Inverse Problem with Denoising Diffusion model priors

Eric Moulines
CMAP, Ecole polytechnique

joint work with Alain Durmus, Lisa Bedin, Gabriel Cardoso, Yazid Janati, Badr Moufad (CMAP), Sylvain Le Corff (LPSM), Jimmy Olsson (KTH)

Introduction

Generative modeling

We have a dataset $\mathcal{D}_N := \{X^1, \dots, X^N\}$, where $X^i \in \mathbb{R}^{d_x}$.

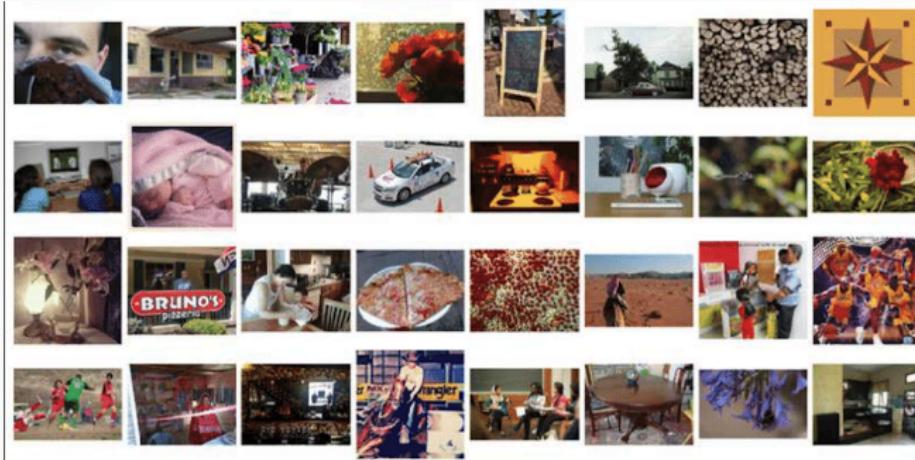


Figure 1: Samples from the ImageNet dataset.

Generative modeling

We have a dataset $\mathcal{D}_N := \{X^1, \dots, X^N\}$, where $X^i \in \mathbb{R}^{d_x}$.

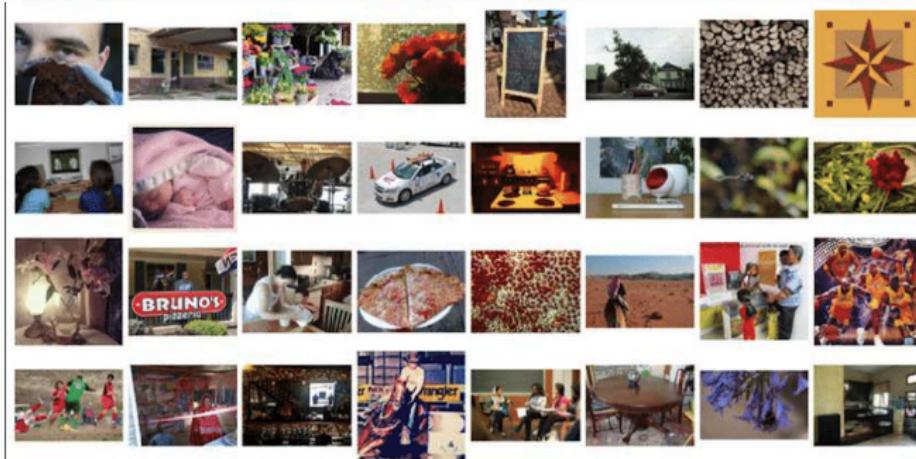


Figure 1: Samples from the ImageNet dataset.

Modeling assumption

(X^1, \dots, X^N) are samples from some **unknown** distribution π_{data}

Generative modeling

- ① Approximate π_{data} with a parametric model.

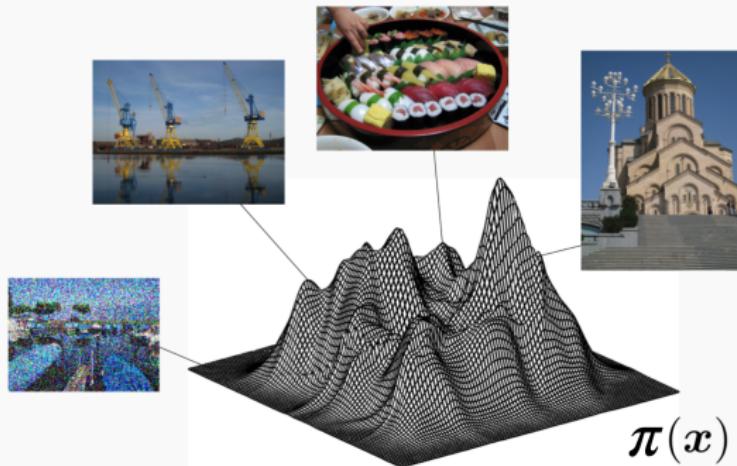


Figure 2: **data** distribution.

Bayesian inverse problems

② Sample reconstructions from the posterior distribution.



Figure 3: Reconstruction problems. Figure adapted from [Lugmayr et al. \(2022\)](#).

Generative modeling

- ① Approximate π_{data} with a **parametric** model p^θ .

Ackley et al. (1985); Kingma and Welling (2013); Goodfellow et al. (2014); Rezende and Mohamed (2015); Sohl-Dickstein et al. (2015); Ho et al. (2020); Song and Ermon (2020)

Generative modeling

① Approximate π_{data} with a **parametric** model p^θ .

1 Choose a **suitable parametric form** for p^θ .

Ackley et al. (1985); Kingma and Welling (2013); Goodfellow et al. (2014); Rezende and Mohamed (2015); Sohl-Dickstein et al. (2015); Ho et al. (2020); Song and Ermon (2020)

Generative modeling

① Approximate π_{data} with a **parametric** model p^θ .

- 1 Choose a **suitable parametric form** for p^θ .
- 2 Train p^θ to approximate π using the samples $(X^1, \dots, X^N) \sim \pi$.

$$\mathcal{L}(\theta) = \sum_{i=1}^N -\log p^\theta(X^i).$$

\rightsquigarrow Minimize $\mathcal{L}(\theta) \rightarrow$ find optimal parameter θ_* .

Ackley et al. (1985); Kingma and Welling (2013); Goodfellow et al. (2014); Rezende and Mohamed (2015); Sohl-Dickstein et al. (2015); Ho et al. (2020); Song and Ermon (2020)

Posterior sampling

② Perform controlled generation using p^{θ_*} .

~~ Target distribution: weight p^{θ_*} with a function $x \mapsto g(x)$

Posterior sampling

② Perform controlled generation using p^{θ_*} .

~~ Target distribution: weight p^{θ_*} with a function $x \mapsto g(x)$

$$\phi(dx) = \frac{g(x)p^{\theta_*}(dx)}{\int g(z)p^{\theta_*}(dz)},$$

~~ Posterior sampling: $g(x) = p(y|x)$.

~~ Reinforcement learning: g is a reward function.

Denoising diffusion models

Introduction

- A denoising diffusion probabilistic model (DDPM) makes use of two Markov chains:
 - 1 a forward chain (process) that perturbs data to noise,
 - 2 a reverse chain (process) that converts noise back to data.
- The forward chain is typically hand-designed with the goal to transform the data distribution π_{data} into a (simple) reference distribution π_{ref} (e.g., standard Gaussian)
- The backward chain reverses the forward chain by learning transition kernels.
- New data points are generated by first sampling a random vector from the reference distribution, followed by ancestral sampling through the backward Markov chain.

Forward process

- Given a data distribution $x_0 \sim \pi_{\text{data}}(dx_0) = q_0(dx_0)$, the **forward Markov chain** generates a sequence of random variables $x_1, x_2 \dots x_T$ with transition kernel $q_{t|t-1}(dx_t | x_{t-1})$.
- The joint distribution of $x_1, x_2 \dots x_T$ conditioned on x_0 , denoted as $q_{0:T}(d(x_1, \dots, x_T) | x_0)$, may be written as

$$q_{0:T}(d(x_1, \dots, x_T) | x_0) = \prod_{t=1}^T q_{t|t-1}(dx_t | x_{t-1}).$$

- In DDPMs, we handcraft the transition kernel $q_{t|t-1}(dx_t | x_{t-1})$ to incrementally transform the data distribution $q_0(dx_0)$ into a tractable **reference distribution**.

- Typical design: Gaussian perturbation

$$q_{t|t-1}(x_t | x_{t-1}) = \mathcal{N}\left(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}\right),$$

where $\beta_t \in (0, 1)$ is a hyperparameter chosen ahead of model training.

Forward process

- Gaussian transition kernel allows us to obtain the analytical form of $q_{t|0}(x_t | x_0)$ for all $t \in \{0, 1, \dots, T\}$. Setting $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=0}^t \alpha_s$, we have

$$q_{t|0}(x_t | x_0) = \mathcal{N}\left(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}\right).$$

- Given x_0 , we can easily obtain a sample of x_t by sampling a Gaussian vector $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$ and applying the transformation

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t.$$

- When $\bar{\alpha}_T \approx 0$, x_T is almost Gaussian in distribution,

$$q_T(x_T) := \int q_{T|0}(x_T | x_0) q_0(x_0) dx_0 \approx \mathcal{N}(x_T; \mathbf{0}, \mathbf{I}).$$

Backward process

- For generating new data samples, DDPMs start by sampling the **reference distribution** and then **gradually remove noise** by running a **learnable Markov chain backward in time**.
- The reverse Markov chain is parameterized by a **reference distribution** $\pi_{\text{ref}}(x_T) = \mathcal{N}(x_T; \mathbf{0}, \mathbf{I})$ and a **learnable transition kernel**

$$p_{t-1|t}^\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \mu_t^\theta(x_t), \Sigma_t^\theta(x_t))$$

where θ denotes model parameters, and the mean $\mu_t^\theta(x_t)$ and variance $\Sigma_t^\theta(x_t)$ are parameterized by deep neural networks.

■ Data generation

- Sample $x_T \sim \pi_{\text{ref}}(\cdot)$,
- iteratively sample $x_{t-1} \sim p_{t-1|t}^\theta(\cdot | x_t)$ until $t = 1$.

Diffusion model principles

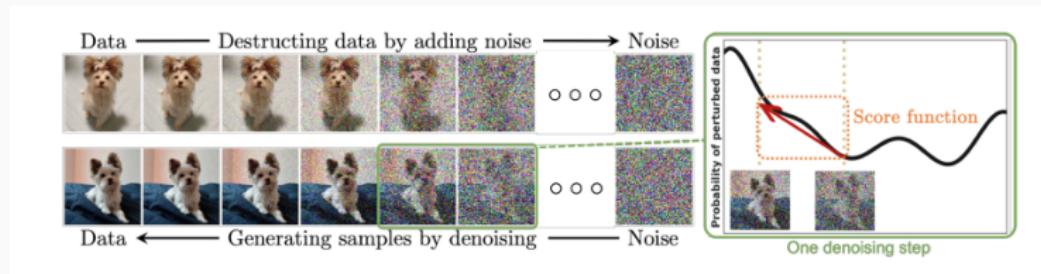


Figure 4: Diffusion models smoothly perturb data by adding noise, then reverse this process to generate new data from noise.

Variational Inference

- **Objective:** Adjust the parameter θ so that the joint distribution of the reverse Markov chain

$$p_{0:T}^{\theta}(x_0, x_1, \dots, x_T) = p_{\text{ref}}(x_T) \prod_{t=1}^T p_{t-1|t}^{\theta}(x_{t-1} | x_t)$$

matches

$$q_{0:T}(x_0, x_1, \dots, x_T) := q_0(x_0) \prod_{t=1}^T q_{t|t-1}(x_t | x_{t-1}).$$

- Training is performed by maximizing a **variational bound**:

$$\begin{aligned} \mathbb{E}_{q_0}[-\log p^{\theta}(x_0)] &\leq \mathbb{E}_{q_{0:T}}\left[-\log \frac{p_{0:T}^{\theta}(x_{0:T})}{q_{1:T|0}(x_{1:T} | x_0)}\right] \\ &= \mathbb{E}_{q_{0:T}}\left[-\log p_T(x_T) - \sum_{t \geq 1} \log \frac{p_{t-1|t}^{\theta}(x_{t-1} | x_t)}{q_{t|t-1}(x_t | x_{t-1})}\right] =: L^{\theta} \end{aligned}$$

Variational inference with variance reduction

- L^θ might be rewritten using the **backward** representation of the **forward** noising process

$$\begin{aligned} q_{1:T|0}(x_{1:T}|x_0) &= \prod_{t=1}^T q_{t|t-1}(x_t|x_{t-1}) \\ &= q_{T|0}(x_T|x_0) \prod_{t=2}^T q_{t-1|t}(x_{t-1}|x_t, x_0) \end{aligned}$$

- With this backward decomposition, L^θ writes

$$\begin{aligned} L^\theta &= \mathbb{E}_{q_{0:T}} \left[-\log \frac{p_T(x_T)}{q_{T|0}(x_T | x_0)} - \sum_{t=2}^T \log \frac{p_{t-1|t}^\theta(x_{t-1} | x_t)}{q_{t-1|t,0}(x_{t-1} | x_t, x_0)} \right. \\ &\quad \left. - \log p_{0|1}^\theta(x_0 | x_1) \right] \\ &= \mathbb{E}_{q_{0:T}} [D_{\text{KL}}(q_{T|0}(\cdot | x_0) \| p_T(\cdot))] \\ &\quad + \sum_{t=2}^T D_{\text{KL}} \left(q_{t-1|t,0}(\cdot | x_t, x_0) \| p_{t-1|t}^\theta(\cdot | x_t) \right) - \log p_{0|1}^\theta(x_0 | x_1) \end{aligned}$$

Variational inference with variance reduction

- forward posteriors are tractable when conditioned on x_0 :

$$q_{t-1|t,0}(x_{t-1} | x_t, x_0) = \mathcal{N}\left(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t \mathbf{I}\right)$$

$$\text{where } \tilde{\mu}_t(x_t, x_0) := \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}x_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}x_t$$

$$\text{and } \tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t$$

- KL divergences are comparisons between Gaussian distributions with closed form expressions: taking $\Sigma_t^\theta(x_t) = \tilde{\beta}_t \mathbf{I}$,

$$D_{\text{KL}}\left(q_{t-1|t,0}(\cdot | x_t, x_0) \| p_{t-1|t}^\theta(\cdot | x_t)\right) = \frac{1}{2\tilde{\beta}_t} \|\tilde{\mu}_t(x_t, x_0) - \mu_t^\theta(x_t)\|^2.$$

Variational inference with variance reduction

- Setting

$$\mu_t^\theta(x_t) = \tilde{\mu}_t(x_t, \hat{x}_{0|t}^\theta(x_t)),$$

we get

$$D_{\text{KL}} \left(q_{t-1|t,0}(\cdot | x_t, x_0) \| p_{t-1|t}^\theta(\cdot | x_t) \right) = w_t \|x_0 - \hat{x}_{0|t}^\theta(x_t)\|^2.$$

with $w_t = \bar{\alpha}_{t-1}\beta_t/(1 - \bar{\alpha}_{t-1})(1 - \bar{\alpha}_t)$.

- Hence, criterion L^θ rewrites

$$L^\theta = \sum_{t=2}^T w_t \mathbb{E}_{q_0 \otimes \mathcal{N}(0, I)} [\|x_0 - \hat{x}_{0|t}^\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon)\|^2]$$

which amounts to compute $\hat{x}_{0|t}^\theta(x_t)$ as a **predictor** of the initial state x_0 from the current state x_t .

- This criterion is the **denoising score matching**.

Noise prediction

- Using that $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t$, we have

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(x_t - \sqrt{1 - \bar{\alpha}_t}\epsilon_t)$$

- Choosing $\hat{x}_{0|t}^\theta(x_t) = (1/\sqrt{\bar{\alpha}_t})(x_t - \sqrt{1 - \bar{\alpha}_t}\hat{\epsilon}_{0|t}^\theta(x_t))$, the criterion L^θ may be equivalently expressed as

$$L^\theta = \sum_{t=2}^T \tilde{w}_t \mathbb{E}_{q_0 \otimes \mathcal{N}(0, I)} [\|\epsilon - \hat{\epsilon}_{0|t}^\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon)\|^2]$$

where

$$\tilde{w}_t = \frac{\beta_t}{\alpha_t(1 - \bar{\alpha}_{t-1})}$$

A continuous-time perspective

Ornstein-Uhlenbeck Noising process

- Consider a diffusion process $\{X_t\}_{t=0}^T$ that starts from the data distribution $q_0(dx) \equiv \pi_{\text{data}}(dx)$ at time $t = 0$. The notation $q_t(dx)$ refers to the marginal distribution of the diffusion at time $0 \leq t \leq T$.
- Assume furthermore that at time $t = T$, the marginal distribution is (very close to) a reference distribution $q_T(dx) = \pi_{\text{ref}}(dx)$ that is straightforward to sample from, e.g. $\mathcal{N}(0, I)$.
- This diffusion process is the **noising process**. It is often chosen as an **Ornstein-Uhlenbeck (OU)** diffusion,

$$dX_t = -\frac{1}{2}X_t dt + dW_t$$

OU noising process

- OU diffusion is **reversible** w.r.t. $\pi_{\text{ref}} = \mathcal{N}(0, I)$: the conditional distribution of $X_{t+s} | X_t = x_t$ is $\mathcal{N}(\alpha_s x_t, \sigma_s^2 I)$, with

$$\alpha_s = \sqrt{1 - \sigma_s^2} \quad \sigma_s^2 = 1 - e^{-s}$$

- Denote

$$F(s, x, y) \propto \exp \left\{ -\frac{(y - \alpha_s x)^2}{2\sigma_s^2} \right\}.$$

the **forward transition** from x to y in " s " amount of time.

Reverse diffusion I (informal)

- the DDPM strategy consists in sampling from the Gaussian reference measure π_{ref} at time $t = T$ and simulate the OU process backward in time.
- In other words, one would like to simulate from the reverse process \overleftarrow{X}_t defined as

$$\overleftarrow{X}_s = X_{T-s}$$

- The reverse process is distributed as $\overleftarrow{X}_0 \sim \pi_{\text{ref}}$ at time $t = 0$ and, crucially, we have that $\overleftarrow{X}_T \sim \pi_{\text{data}}$.
- The reverse diffusion follows the dynamics (Hausmann, Pardoux, 1986; Millet, Nualart, Sanz, 1989)

$$d\overleftarrow{X}_t = +\frac{1}{2}\overleftarrow{X}_t dt + \nabla \log q_{T-t} \left(\overleftarrow{X}_t \right) dt + dB_t$$

where B is another Wiener process [the notation B emphasizes that there is no link between this Wiener process and the one used to simulate the forward process]

Reverse diffusion II (informal)

- To simulate the reverse diffusion, one needs to be able to estimate the **score** $\nabla \log q_{T-t}(x)$.
- In practice, the score is **unknown** and need to be **approximated**

$$s_t^\theta(x) \approx \nabla_x \log q_t(x)$$

which is often parameterized by a neural network.

- Since

$$\log q_t(x) = \log \int F(t, x_0, x) \pi_{\text{data}}(dx_0)$$

the analytical expression of $F(t, x_0, x)$ gives that (**Tweedie formula**)

$$\nabla_x \log q_t(x) = -\frac{x - \alpha_t \hat{x}_{0|t}(x)}{\sigma_t^2}$$

where $\hat{x}_{0|t}(x) = \mathbb{E}[X_0 | X_t = x]$ is a **denoising** estimate of x_0 given a **noisy estimate** $X_t = x$ at time t

Estimation of the score

- To estimate the score, one only needs to train a denoising function $\hat{x}_{0|t}^\theta(x)$.
- It is a simple regression problem: take pairs (X_0, X_t) that can be generated as

$$X_0 \sim \pi_{\text{data}} \quad \text{and} \quad X_t = \alpha_t X_0 + \sigma_t Z_t$$

with $Z_t \sim \mathcal{N}(0, I)$ and minimize the Mean Squared Error (MSE) loss, i.e.

$$\mathbb{E}_{q_{0,t}} \left[\left\| X_0 - \hat{x}_{0|t}^\theta(X_t) \right\|^2 \right]$$

with stochastic gradient descent or any other stochastic optimization procedure.

- The score is then defined as

$$s_t^\theta(x) = -\frac{x - \alpha_t \hat{x}_{0|t}^\theta(x)}{\sigma_t^2}$$

Time reversal formula for a diffusion process

General time reversal formulas for diffusion processes are well known since the 80 's. Consider a diffusion process Y in \mathbb{R}^n satisfying

$$dY_t = b_t(Y_t) dt + \sigma_t(Y_t) dB_t, \quad 0 \leq t \leq T,$$

with B a Brownian motion, b a drift vector field and σ a matrix field associated to the diffusion field $a := \sigma\sigma^t$, (σ^t is the transposed of σ .)

Assuming that the law of Y_t is absolutely continuous at each time t , under various hypotheses on b and a , one can prove that the time-reversed process Y^* is again a diffusion process with diffusion matrix field $a_t^* = a_{T-t}$ and drift field

$$b_t^*(y) = -b_{T-t}(y) + \nabla \cdot (\mu_{T-t} a_{T-t})(y) / \mu_{T-t}(y),$$

where μ_t is the density of the law of Y_t with respect to Lebesgue measure. This is not a straightforward result because a reversed semimartingale might not be a semimartingale anymore.

Time reversal formula for a diffusion process

Assumptions

- b is locally Lipschitz,
- a is bounded away from zero or that the derivative ∇a is well-behaved.

then

$$b_t^*(y) = -b_{T-t}(y) + \nabla \cdot (\mu_{T-t} a_{T-t})(y) / \mu_{T-t}(y),$$

- Haussmann and Pardoux take a PDE approach;
- Millet, Nualart and Sanz rely on stochastic calculus of variations.
- The existence of an absolutely continuous density follows from a Hörmander type condition (PDE formulation in Haussman et al. and consequence of Malliavin calculus in Millet et al.).

Summary

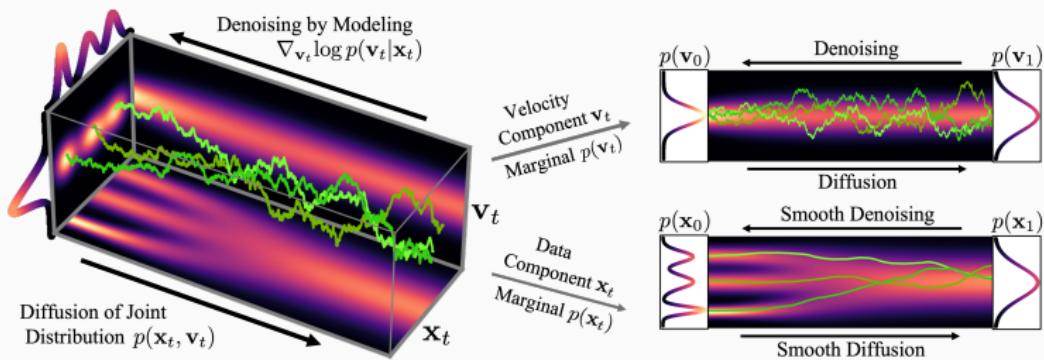


Figure 5: From Dockhorn et al. (2022)

Posterior Sampling

Context

Bayesian linear inverse problem:

$$Y = AX + \sigma_y Z, \quad \text{where} \quad Z \sim \mathcal{N}(\mathbf{0}_{d_x}, \mathbf{I}_{d_x}), \quad X \sim p_0, \quad \sigma_y \geq 0.$$

Context

Bayesian linear inverse problem:

$$Y = AX + \sigma_y Z, \quad \text{where} \quad Z \sim \mathcal{N}(\mathbf{0}_{d_x}, \mathbf{I}_{d_x}), \quad X \sim p_0, \quad \sigma_y \geq 0.$$

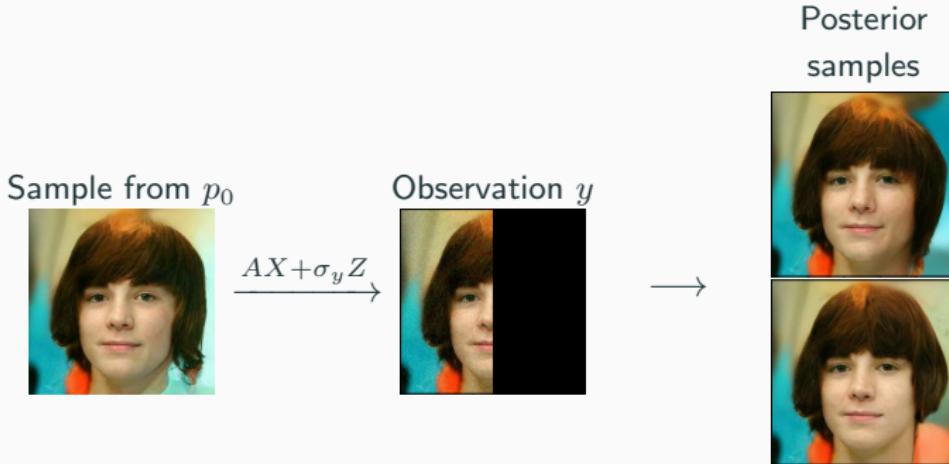
Objective: Sample the distribution of X given a realisation y of Y .

Context

Bayesian linear inverse problem:

$$Y = AX + \sigma_y Z, \quad \text{where} \quad Z \sim \mathcal{N}(\mathbf{0}_{d_x}, \mathbf{I}_{d_x}), \quad X \sim p_0, \quad \sigma_y \geq 0.$$

Objective: Sample the distribution of X given a realisation y of Y .



Inverse problems

Goal: Reconstruct an unknown \mathbb{R}^{d_x} -valued signal X from a \mathbb{R}^m -valued measurement Y based on a corruption model given by

$$Y = \mathcal{A}(X) + \sigma Z, \quad Z \sim \mathcal{N}(0, I_m),$$

where \mathcal{A} is the forward operator, Z is unobserved noise, and $\sigma > 0$ denotes the noise level.

- **denoising**, where $\mathcal{A}(x) = x$ and $\sigma > 0$ and the noise does not need to be Gaussian but can be heavy-tailed or even multiplicative;
- **inpainting**, where $\mathcal{A}(x) = Ax$ where A is a masking matrix with $A_{ij} = 0$ for $i \neq j$ and A_{ii} is either 0 or 1,
- **deblurring**, where $\mathcal{A}(x) = Ax$ where A is a 2D convolution operator (point spread function);
- **phase retrieval**, where $\mathcal{A}(x) = |Ax|$, where A is an invertible matrix.

III-posedness

- A key feature of these problems is that they are **ill-posed**, which means that some information is lost, making exact recovery impossible.
- In the Bayesian inverse problem framework, the solution is treated as a **distribution**. This approach represents unknown parameters as a random variables and updates the beliefs using Bayes' theorem based on the observed data.
- The Bayesian framework enables the incorporation of prior knowledge and the quantification of uncertainty, which is essential for addressing ill-posed problems
 - There is generally not a single solution compatible with the observations, but rather a set of solutions whose plausibility is measured by the posterior distribution.

Bayesian inverse problems

- In the Bayesian approach to inverse problems, we provide the unknown signal X with a prior distribution q .
- The likelihood function is expressed as

$$p(y|x) \propto \exp(-\|y - \mathcal{A}(x)\|^2/2\sigma^2).$$

- The goal is to sample from the posterior distribution
 $p(x|y) \propto p(y|x)q(x)$.
- This survey focuses on a specific class of methods that employ denoising diffusion models (DDMs) as priors. The goal is to sample reconstructions from the posterior $p(x|y)$ using models pre-trained with data from a given distribution, without the need for retraining or fine-tuning [Song and Ermon \(2019\)](#); [Ho et al. \(2020\)](#); [Song et al. \(2021b,a\)](#).

Bayesian inverse problems

- The setup we consider makes it possible, given a pre-trained DDM, to solve an arbitrary Bayesian inverse problem using this DDM as a prior, which eliminates the need to train a task-specific conditional model from scratch.
 - This is a departure from typical conditional DDM frameworks [Song et al. \(2021b\)](#); [Batzolis et al. \(2021\)](#); [Tashiro et al. \(2021\)](#); [Saharia et al. \(2022\)](#), which require access to a paired data set (X, Y) to learn a parametric approximation of the posterior distribution of the signal X given the observation Y .

Bayesian inverse problems

In the sequel we consider the problem of sampling approximately from some target distribution

$$\pi(dx) := g(x)p(dx)/\mathcal{Z}$$

on \mathbb{R}^{d_x} , where

- g is some non-negative function on \mathbb{R}^{d_x} referred to as the **potential**, which is assumed to be evaluable pointwise,
- p is a prior distribution, and $\mathcal{Z} := \int g(x)p(dx)$ is the normalising constant. We assume that p is provided by a DDM.

DDM prior

- A DDM prior imposes constraints that limit the flexibility in designing a sampler for the target posterior. First, recall that the posterior distribution π , depends on the prior p , which is intractable, and we only have access to its parametric approximation p^θ .
- Since p_0^θ represents the marginal of the latent variable model, we cannot directly evaluate its density or reliably approximate its score $\nabla \log p_0^\theta$.
- This makes it impractical to bridge between a standard Gaussian distribution and π via an annealing path for example; i.e., by considering an interpolation of the form $\pi_k \propto \pi^{\gamma_k} \mathcal{N}(\cdot; \mathbf{0}_{d_x}, \mathbf{I}_{d_x})^{1-\gamma_k}$, where $(\gamma_k)_{k=n}^0$ is a temperature schedule satisfying $\gamma_n = 0$ and $\gamma_0 = 1$ [Syed et al. \(2021\)](#); [Dai et al. \(2022\)](#).

DDM prior

- A natural approach to constructing a sampler for the posterior distribution π is to leverage the pretrained backward kernels $(p_{k|k+1}^\theta(x_k|x_{k+1}))_{k=0}^n$.
- We define the intermediate posterior distribution $\pi_k^*(x_k)$ as follows:

$$\pi_k^*(x_k) := \int q_{k|0}(x_k|x_0) \pi(dx_0) = \int \frac{g(x_0)p(dx_0)}{p(g)} q_{k|0}(x_k|x_0).$$

This is the analogue of forward marginals, with the crucial difference being the replacement of p by π .

- However, unlike the marginal distribution of the forward process, $q_k(dx_k)$, sampling from this posterior distribution is intractable, as direct sampling from π is not possible.

Posterior sequence

- Idea: relate $\pi_k^*(x_k)$ to $q_k(dx_k)$ through the following expression:

$$\pi_k^*(x_k) = \frac{g_k(x_k)q_k(x_k)}{q_k(g_k)}, \quad \text{where} \quad g_k(x_k) = \int g(x_0)q_{0|k}(dx_0|x_k).$$

where we have set

$$q_{0|k}(dx_0|x_k) := p(dx_0)q_{k|0}(dx_k|x_0)/q_k(x_k).$$

- The structure of π_k^* closely mirrors that of the posterior distribution π . Both are expressed as the product of the forward process's k -th marginal, q_k , and a potential function, which can be interpreted as the likelihood of the observation given the forward process at time k .

Conditional score and guidance terms

- The conditional scores $(\nabla \log \pi_k^*)_{k=1}^n$ can be obtained similarly with q_k replaced by

$$\nabla_{x_k} \log \pi_k^*(x_k) = \nabla_{x_k} \log q_k(x_k) + \nabla_{x_k} \log g_k(x_k).$$

- From Tweedie's formula, the conditional score is also related to the conditional denoiser $\mathbb{E}^\pi[X_0|X_k = x_k] = \int x_0 \pi_{0|k}^*(dx_0|x_k)$ where $\pi_{0|k}^*(dx_0|x_k) \propto \pi(dx_0)q_{k|0}(x_k|x_0)$.
 - The term $\nabla_{x_k} \log g_k(x_k)$, often referred to in the literature as the **guidance term**, indicates how the original score should be adjusted to steer the samples towards the posterior distribution.

Gradient guidance

- **Idea:** Approximate the denoiser $(x, k) \mapsto \mathbb{E}^\pi[X_0 | X_k = x]$ where the conditional expectation is taken with respect to $\pi_{0|k}^*(dx_0|x_k) \propto \pi(dx_0)q_{k|0}(x_k|x_0)$.

- **Key relation:**

$$\mathbb{E}^\pi[X_0 | X_k = x_k] = \mathbb{E}^p[X_0 | X_k = x_k] + \frac{1 - \alpha_k}{\sqrt{\alpha_k}} \nabla_{x_k} \log g_k(x_k)$$

- Using an (unconditional) pre-trained denoisers $(x_{0|k}^\theta)_{k=1}^n$, we can approximate the denoiser $\mathbb{E}^\pi[X_0 | X_k]$ and enable approximate sampling from π using a DDPM-like approximation.

Gradient guidance

- A sample from the DDPM approximation of the posterior transition $\pi_{k|k+1}^*(\cdot|X_{k+1})$ is given by:

$$X_k = \tilde{X}_k + \frac{v_{k+1|k}}{\sqrt{1 - v_{k+1|k}}} \widehat{\nabla \log g}_{k+1}(X_{k+1}), \quad \tilde{X}_k \sim p_{k|k+1}^\theta(\cdot|X_{k+1}).$$

- Most of the methods mimick this update take the form:

$$X_k = \tilde{X}_{k+1} + \omega_k(X_{k+1}) \widehat{\nabla \log g}_{k+1}(X_{k+1}),$$

where ω_k is a weight function that varies according to the specific method.

Diffusion Posterior Sampling

- Chung et al. (2023) suggested to set

$$\nabla_{x_k} \log g_k(x_k) \approx \nabla_{x_k} \log g(x_{0|k}^\theta(x_k)),$$

which amounts to replace $q_{0|k}((|x_k)\mathrm{d}x_0)$ by $\delta_{x_{0|k}^\theta(x_k)}(\mathrm{d}x_0)$ in the definition the posterior.

- It involves the computation of a vector-Jacobian product involving the denoiser network since

$$\nabla_{x_k} \log g(x_{0|k}^\theta(x_k)) = \nabla_{x_k} x_{0|k}^\theta(x_k)^\top \nabla_{x_0} \log g(x_0)|_{x_0=x_{0|k}^\theta(x_k)}.$$

- With $g(x) = \mathcal{N}_{d_y}(y; \mathcal{A}(x), \sigma_y^2 \mathbf{I})$, the update is

$$\omega_k(x_{k+1}) = \gamma \sigma_y^2 / \|y - \mathcal{A}(x_{0|k+1}^\theta(x_{k+1}))\|, \quad \gamma \in [0, 1],$$

Gaussian approximation

- Song et al. (2023) proposed a Gaussian approximation:

$$q_{0|k}(x_0|x_k) \approx \mathcal{N}_{d_x}(x_0; x_{0|k}^\theta(x_k), v_{0|k}\mathbf{I}),$$

where the variance $v_{0|k}$ is treated as a hyperparameter.

- Song et al. (2023) recommend setting $v_{0|k} = 1 - \alpha_k$, which matches the variance of $q_{0|k}(x_0|x_k)$ under the assumption that p is a standard Gaussian distribution.
 - For the specific case of linear inverse problems, where $g(x) = \mathcal{N}_{d_y}(y; Ax, \sigma_y^2\mathbf{I})$, the following approximation is employed.

$$\nabla_{x_k} \log g_k(x_k) \approx \nabla_{x_k} \log \mathcal{N}_{d_y}(y; Ax_{0|k}^\theta(x_k), v_{0|k}AA^\top + \sigma_y^2\mathbf{I})$$

which is obtained by integrating g against the Gaussian approximation (??) using (Bishop, 2006, Formula 2.115).

Improved Gaussian approximation

- The vanilla Gaussian approximation can be further improved: note that the covariance is not related to the data distribution.
- In [Finzi et al. \(2023\)](#); [Boys et al. \(2023\)](#), it is suggested to solve the following optimization problem:

$$\operatorname{argmin}_{\mu \in \mathbb{R}^{d_x}, \Sigma \in \mathcal{S}_{++}^{d_x}} \text{KL}(\mathcal{N}_{d_x}(\cdot; \mu, \Sigma) \parallel q_{0|k}(\cdot | x_k)),$$

where $\mathcal{S}_{++}^{d_x}$ is the set of positive definite matrices.

- Using ([Bishop, 2006](#), Section 10.7), the optimal solution, known as *moment projection*, is given by

$$\mu_{0|k}(x_k) = \mathbb{E}[X_0 | X_k = x_k], \quad \Sigma_{0|k}(x_k) = \text{Cov}[X_0 | X_k = x_k].$$

Feynman-Kac representation

- Although the sequence $(\pi_k^*)_{k=n}^0$ of distributions is a natural choice, this approach poses numerical challenges, as it involves quantities that are difficult to approximate.
- Alternative sequences can be explored, and as long as there is a way to transition effectively from the $(k + 1)$ th distribution to the k th one, accurate samples from the final distribution π can still be obtained.
- This is a well-proven strategy that has been extensively explored; see, e.g., [Gelman and Meng \(1998\)](#); [Del Moral \(2013\)](#); [Dai et al. \(2022\)](#) .

Feynman-Kac representation

Given a generic sequence $(\pi_k)_{k=n}^0$ of distributions, a common approach involves establishing a recursive relationship between successive distributions π_{k+1} and π_k , such as

$$\pi_k(dx_k) = \frac{1}{Z_k} \int m_{k|k+1}(dx_k | x_{k+1}) w_k(x_k, x_{k+1}) \pi_{k+1}(dx_{k+1}),$$

where

- w_k is a non-negative weight function,
- $Z_k = \int m_{k|k+1}(dx_k | x_{k+1}) w_k(x_k, x_{k+1}) \pi_{k+1}(dx_{k+1})$ is a normalising constant,
- $m_{k|k+1}(\cdot | x_{k+1})$ is a Markov kernel.

Sequential Monte Carlo I

- There are various approaches for sampling from a sequence of distributions satisfying

$$\pi_k(dx_k) = \frac{1}{\mathcal{Z}_k} \int m_{k|k+1}(dx_k \mid x_{k+1}) w_k(x_k, x_{k+1}) \pi_{k+1}(dx_{k+1}),$$

- The most straightforward approach is to **sequential Monte Carlo** (SMC) methods (a.k.a. **particle filters**) is to approximate each distribution π_k in the sequence of interest by a weighted sample $(\omega_k^i, \xi_k^i)_{i=1}^N$ of random draws (the ξ_k^i 's), referred to as particles, and associated non-negative importance weights (the ω_k^i 's).

Sequential Monte Carlo II

- An SMC algorithm recursively propagates the weighted particle sample from one time step k to the next using **importance sampling and resampling techniques**, providing an approximation of π_k at each step.
- The particle samples serve two key purposes:
 - first, they can be used to approximate integrals with respect to the target distributions, which allows for estimates of relevant statistics;
 - second, since each particle represents a possible state of the Feynman–Kac flow, the weighted particle samples represent the distributions themselves, enabling inference at subsequent time steps.We refer to [Chopin et al. \(2020\)](#); [Del Moral \(2004\)](#); [Del Moral \(2013\)](#).

SMC for Feynman-Kac I

- Assume that we have access to a sample $(\omega_{k+1}^i, \xi_{k+1}^i)_{i=1}^N$ of particles and associated weights, the weighted empirical distribution π_{k+1}^N of which approximates π_{k+1} .
- In order to form an updated sample $(\omega_k^i, \xi_k^i)_{i=1}^N$ approximating π_k , we simply plug π_{k+1}^N into the FK recursion

$$\pi_k(dx_k) = \frac{1}{Z_k} \int m_{k|k+1}(dx_k | x_{k+1}) w_k(x_k, x_{k+1}) \pi_{k+1}(dx_{k+1}),$$

- This yields a mixture approximating π_k .

$$\bar{\pi}_k^N(dx_k) \propto \sum_{i=1}^N \omega_{k+1}^i w_k(x_k, \xi_{k+1}^i) m_k(dx_k | \xi_{k+1}^i)$$

SMC for Feynman-Kac II

- New particles are sampled from $\bar{\pi}_k^N$ using importance sampling based on some mixture proposal

$$\rho_k^N(dx_k) \propto \sum_{i=1}^N \varphi_{k+1}^i \omega_{k+1}^i r_k(dx_k | \xi_{k+1}^i),$$

- r_k is some transition density dominating m_k
- $(\varphi_{k+1}^i)_{i=1}^N$ is a set of positive weights referred to as adjustment multipliers Pitt and Shephard (1999).

SMC for Feynman-Kac III

- (i) Draw indices $(\kappa_k^\ell)_{\ell=1}^N$ (conditionally) independently from the categorical distribution on $\{1, \dots, N\}$ formed by the adjusted particle weights $(\omega_{k+1}^\ell \varphi_{k+1}^\ell)_{\ell=1}^N$;
- (ii) Generate independently, for all sampled mixture indices, new particles $\xi_k^i \sim r_k(dx_k | \xi_{k+1}^{\kappa_k^i})$;
- (iii) Assign each new particle ξ_k^i an updated weight

$$\omega_k^i := w_k(\xi_k^i, \xi_{k+1}^{\kappa_k^i})(dm_k/dr_k)(\xi_k^i, \xi_{k+1}^{\kappa_k^i})/\varphi_k^{\kappa_k^i},$$

where dm_k/dr_k denotes the Radon–Nikodym derivative of m_k w.r.t. r_k (which is well defined).

SMC for Feynman-Kac IV

- The sampling of the mixture indices (Step (i)) in can be equivalently understood as resampling, with replacement, new particles among the old ones in proportion to the adjusted weights.
- The sampling step (ii) serves to randomly jitter, or shake, the particles. We will therefore refer to operations (i–iii) collectively as the **(resample–shake–weight)** scheme (in the literature, operations (i) and (ii) are also often referred to as selection and mutation, respectively).

SMC for Feynman-Kac V

In some situations, the weight function w_k does not depend on x_k ; in that case, we may sample exactly from $\bar{\pi}_k^N$, without the need of importance sampling, by instead

- (i') computing the weights $(w_k(\xi_{k+1}^i) \omega_{k+1}^i)_{i=1}^N$;
- (ii) drawing independently mixture indices $(\kappa_k^\ell)_{\ell=1}^N$ from the categorical distribution on $\{1, \dots, N\}$ induced by the weights in (i');
- (iii) generating new particles $\xi_k^i \sim m_k(dx_k | \xi_{k+1}^{\kappa_{k+1}^i})$.

Appealingly, this procedure, which we refer to as the **(weight–resample–shake)** scheme, will always generate uniformly weighted particles at all iterations.

Feynman-Kac representation

We focus on the specific case where the prior p_0 is the marginal w.r.t. x_0 of Denoising Diffusion Model. The posterior is

$$p_0^y(dx_0) = \frac{1}{\mathcal{Z}^y} \int g_0^y(x_0) \prod_{k=0}^{n-1} p_{k|k+1}(dx_k|x_{k+1}) p_n(dx_n).$$

- The posterior can be interpreted as the marginal of a (time-reversed) Feynman–Kac (FK) model with **non-trivial potential only at $k = 0$!**
- In this work, we twist, **without modifying the law of the FK model**, the backward transitions $p_{k|k+1}$ by **potentials** depending on the observation y ; see e.g. for a similar idea for rare event simulation (see, e.g., Cérou et al., 2012).

"Forward" smoothing decomposition

- Define, for all $k \in \llbracket 0, n \rrbracket$, the **backward functions**

$$\beta_{0|k}^y(x_k) := \int g_0^y(x_0) p_{0|k}(\mathrm{d}x_0|x_k)$$

- The backward functions satisfy the recursion:

$$\beta_{0|k+1}^y(x_{k+1}) = \int \beta_{0|k}^y(x_k) p_{k|k+1}(\mathrm{d}x_k|x_{k+1}).$$

- Define the **forward smoothing kernels** (FSK) for $k \in \llbracket 0, n-1 \rrbracket$

$$\begin{aligned} p_{k|k+1}^y(\mathrm{d}x_k|x_{k+1}) &:= \frac{\beta_{0|k}^y(x_k)}{\beta_{0|k+1}^y(x_{k+1})} p_{k|k+1}(\mathrm{d}x_k|x_{k+1}), \\ &\quad (= \text{Law}(X_k \mid Y = y, X_{k+1} = x_{k+1})). \end{aligned}$$

“Forward” smoothing decomposition

The posterior distribution can be written in terms of forward smoothing kernels

$$p_0^y(dx_0) = \int p_n^y(dx_n) \prod_{k=0}^{n-1} p_{k|k+1}^y(dx_k|x_{k+1}).$$

where

$$p_n^y(dx_n) = \frac{\beta_{0|n}^y(x_n)p_n(dx_n)}{\mathcal{Z}^y}$$

- Most of the recent works to sample from p_0^y use the **forward smoothing decomposition** with different approximation of the **intractable** forward smoothing kernels. Chung et al. (2023); Song et al. (2023); Zhang et al. (2023); Boys et al. (2023); Trippe et al. (2023); Wu et al. (2023).

Monte Carlo guided diffusion

General Feynman–Kac model

Introduce intermediate positive potentials $(g_k^y)_{k=0}^n$, each being a function on \mathbb{R}^{d_x} , and write

$$\begin{aligned} p_0^y(dx_0) &= \frac{1}{\mathcal{Z}^y} \int g_n^y(x_n) p_n(dx_n) \\ &\quad \times \prod_{k=0}^{n-1} \frac{g_k^y(x_k)}{g_{k+1}^y(x_{k+1})} p_{k|k+1}(dx_k|x_{k+1}). \end{aligned}$$

- Because the $g_n^y(x_n) \prod_{k=0}^{n-1} \frac{g_k^y(x_k)}{g_{k+1}^y(x_{k+1})} = g_0^y(x_0)$, the FK is not modified - the potentials are used to render the sampling easier.
- This allows the posterior of interest to be expressed as the time-zero marginal of a **Feynman-Kac** model with
 - initial law p_n ,
 - Markov transition kernels $(p_{k|k+1})_{k=0}^{n-1}$
 - Potentials g_n^y and $(x_k, x_{k+1}) \mapsto g_k^y(x_k)/g_{k+1}^y(x_{k+1})$.

Posterior sampling proposal

Alternatively, the previous decomposition defines a sequence of distributions

$$p_k^y(dx_k) \propto g_k^y(x_k) p_k(dx_k), \quad k \in [0, n],$$

where the posterior of interest is the terminal distribution at $k = 0$.

- If we have a particle approximation of p_{k+1}^y then we can evolve it into a particle approximation of $p_k^y \rightsquigarrow$ **we recursively build an empirical approximation of p_0^y** .
- The choice of potentials $\{g_k^y\}_{k \in [0, n]}$ is crucial; we need to ensure that p_k^y is close enough to p_{k+1}^y so that we can bridge the intermediate distributions efficiently.

Posterior sampling proposal: recursion

Consider the following particle approximation of p_{k+1}^y

$$p_{k+1}^{N,y} = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_{k+1}^i},$$

Recall that $p_k(\mathrm{d}x_k) = \int p_{k|k+1}(\mathrm{d}x_k|x_{k+1})p_{k+1}(\mathrm{d}x_{k+1})$,

Posterior sampling proposal: recursion

Consider the following particle approximation of p_{k+1}^y

$$p_{k+1}^{N,y} = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_{k+1}^i},$$

Recall that $p_k(\mathrm{d}x_k) = \int p_{k|k+1}(\mathrm{d}x_k|x_{k+1})p_{k+1}(\mathrm{d}x_{k+1})$,

$$p_k^y(\mathrm{d}x_k) = \frac{\int \frac{g_k^y(x_k)}{g_{k+1}^y(x_{k+1})} p_{k|k+1}(\mathrm{d}x_k|x_{k+1}) p_{k+1}^y(\mathrm{d}x_{k+1})}{\int \frac{g_k^y(z_k)}{g_{k+1}^y(z_{k+1})} p_{k|k+1}(\mathrm{d}z_k|z_{k+1}) p_{k+1}^y(\mathrm{d}z_{k+1})},$$

Posterior sampling proposal: recursion

Consider the following particle approximation of p_{k+1}^y

$$p_{k+1}^{N,y} = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_{k+1}^i},$$

Recall that $p_k(\mathrm{d}x_k) = \int p_{k|k+1}(\mathrm{d}x_k|x_{k+1})p_{k+1}(\mathrm{d}x_{k+1})$,

$$p_k^y(\mathrm{d}x_k) = \frac{\int \frac{g_k^y(x_k)}{g_{k+1}^y(x_{k+1})} p_{k|k+1}(\mathrm{d}x_k|x_{k+1}) p_{k+1}^y(\mathrm{d}x_{k+1})}{\int \frac{g_k^y(z_k)}{g_{k+1}^y(z_{k+1})} p_{k|k+1}(\mathrm{d}z_k|z_{k+1}) p_{k+1}^y(\mathrm{d}z_{k+1})},$$

and hence

$$p_k^y(\mathrm{d}x_k) \propto \underbrace{\int \frac{g_k^y(z_k)p_k(\mathrm{d}z_k|x_{k+1})}{g_{k+1}^y(x_{k+1})}}_{:=\tilde{\omega}_k(x_{k+1})} p_k^y(\mathrm{d}x_k|x_{k+1}) p_{k+1}^y(\mathrm{d}x_{k+1}),$$

where $p_k^y(\mathrm{d}x_k|x_{k+1}) \propto g_k^y(x_k)p_{k|k+1}(\mathrm{d}x_k|x_{k+1}) \rightarrow$ available in closed form if we use a Gaussian potential with mean linear in x_k .

Posterior sampling proposal: SMC approximation

$$p_k^y(dx_k) = \int p_k^y(dx_k|x_{k+1}) \frac{\tilde{\omega}_k(x_{k+1}) p_{k+1}^y(dx_{k+1})}{\int \tilde{\omega}_k(z_{t+1}) p_{k+1}^y(dz_{k+1})},$$

Assume $p_k^{N,y} = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_{k+1}^i}$ is a particle approximation of $p_{k+1}^{N,y}$.

~~~ **Weight:**

$$p_k^{N,y}(\cdot) \approx \sum_{i=1}^N \frac{\tilde{\omega}_k(\xi_{k+1}^i)}{\sum_{j=1}^N \tilde{\omega}_k(\xi_{k+1}^j)} p_k^y(\cdot|\xi_{k+1}^i).$$

~~~ **Resample:** Draw  $A_{k+1}^{1:N} \stackrel{\text{iid}}{\sim} \text{Categorical}(\{\omega_k^j\}_{j=1}^N)$  where  $\omega_k^j \propto \tilde{\omega}_t(\xi_{k+1}^j)$ .

~~~ **Mutate:** Sample  $\xi_k^i \sim p_k^y(\cdot|\xi_{k+1}^{A_{k+1}^i})$  for  $i \in [1 : N]$ ,

$$p_k^{N,y} = \frac{1}{N} \sum_{i=1}^N \delta_{\xi_k^i}.$$

---

Gordon et al. (1993); Del Moral (2004); Cappe et al. (2005); Chopin et al. (2020)

# Potentials: heuristic

For simplicity (and only in this slide) let  $p_0(y)$  be the posterior of the inverse problem

$$Y = \bar{X}_0, \quad X_0 \sim p_0,$$

The marginals of the *forward process* initialized at  $p_0^y$  are

$$X_k \stackrel{\mathcal{L}}{=} \sqrt{\bar{\alpha}_k} X_0 + \sqrt{1 - \bar{\alpha}_k} Z, \quad \textcolor{brown}{X}_0 \sim \textcolor{brown}{p}_0^y, \quad Z \sim \mathcal{N}(\mathbf{0}_{d_x}, \mathbf{I}_{d_x}),$$

and so

$$\bar{X}_k \stackrel{\mathcal{L}}{=} \sqrt{\bar{\alpha}_k} y + \sqrt{1 - \bar{\alpha}_k} \bar{Z}, \quad \bar{Z} \sim \mathcal{N}(\mathbf{0}_{d_y}, \mathbf{I}_{d_y}).$$

- This suggests that one relevant choice of potentials is

$$g_k^y(x_k) = \mathcal{N}(\textcolor{red}{\sqrt{\bar{\alpha}_k} y}; x_k, (1 - \bar{\alpha}_k) \mathbf{I}_{d_y}).$$

# Choice of potentials

- More generally, we let the variance be a **free parameter**  $\sigma_{y,k}^2$ .

*Our proposal in the general case is*

$$p_k^y(dx_k) \propto g_k^y(x_k) p_k(dx_k), \quad g_k^y(x_k) := \mathcal{N}(\sqrt{\alpha_k} y; Ax_k, \sigma_{y,k}^2 I_{d_y})$$

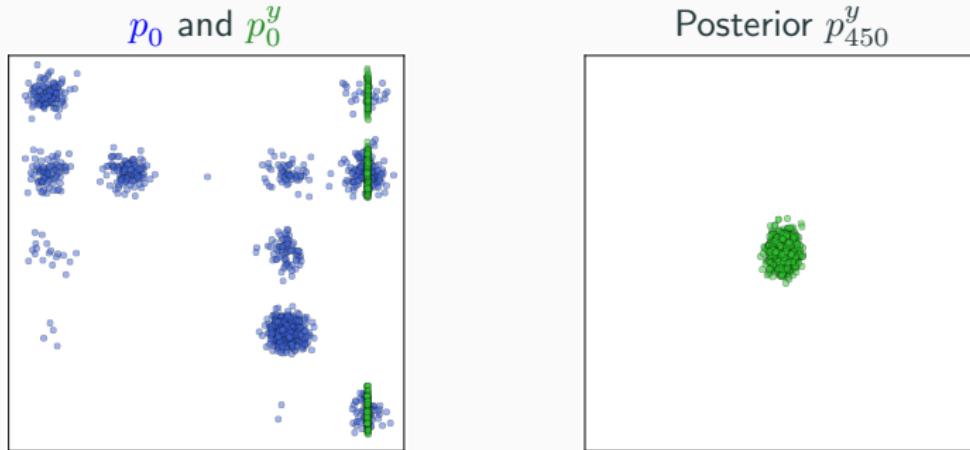
- This particular choice of potential allows us to compute in closed form the auxiliary transition kernel  $\propto g_k^y(x_k) p_{k|k+1}(dx_k|x_{k+1})$  we use for our particle approximations.

# Illustration

$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.

# Illustration

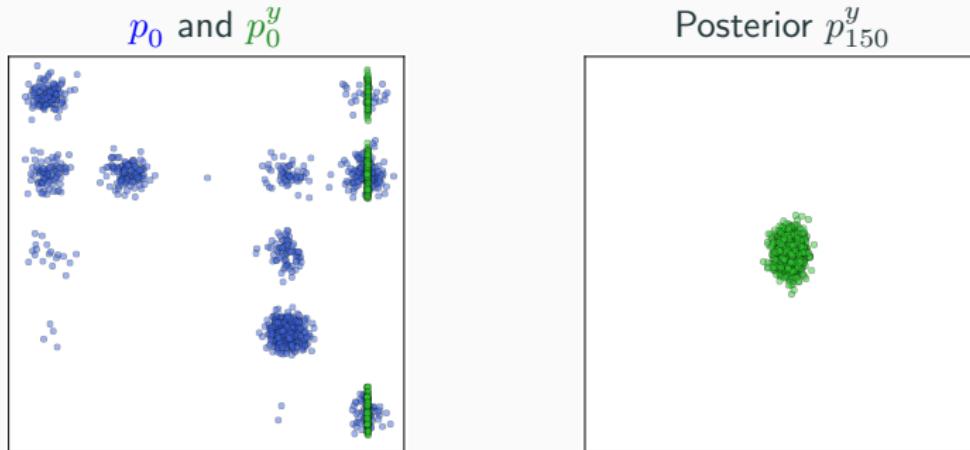
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 6:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

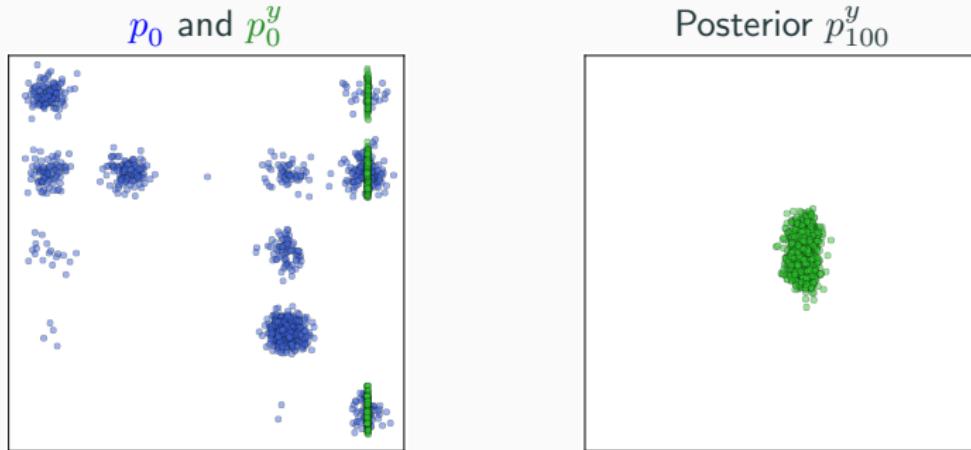
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 7:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

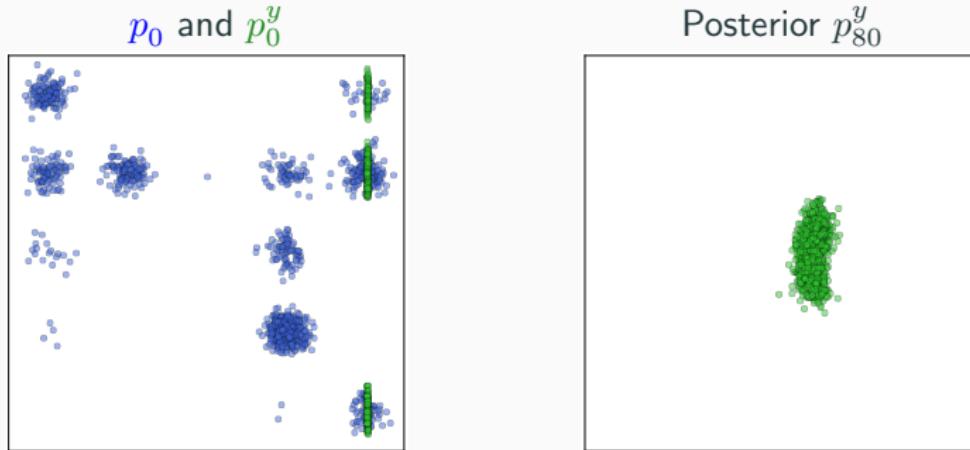
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 8:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

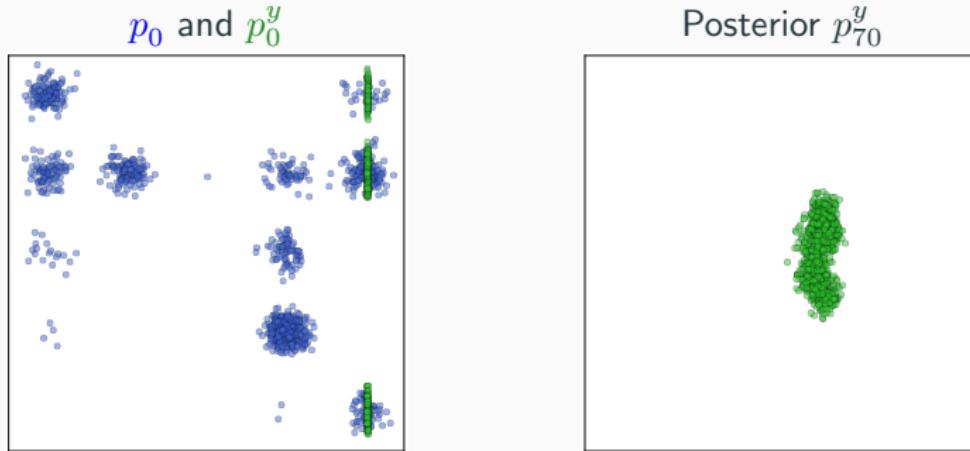
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 9:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

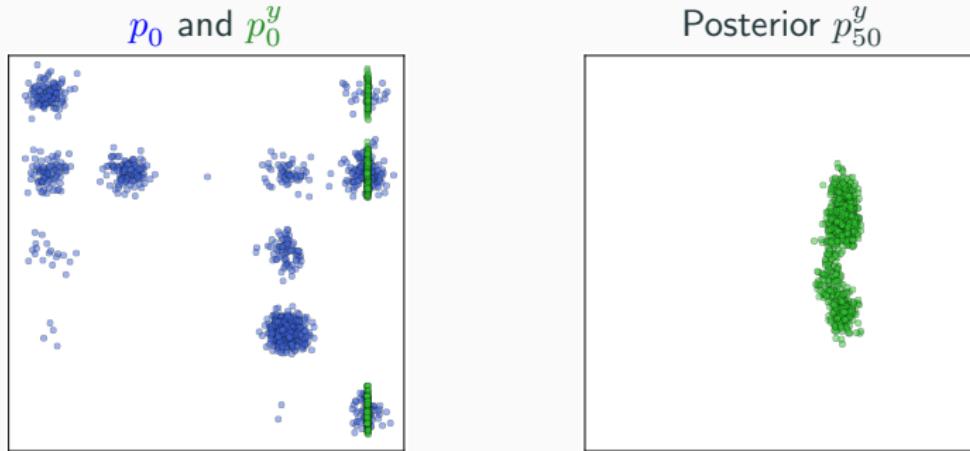
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 10:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

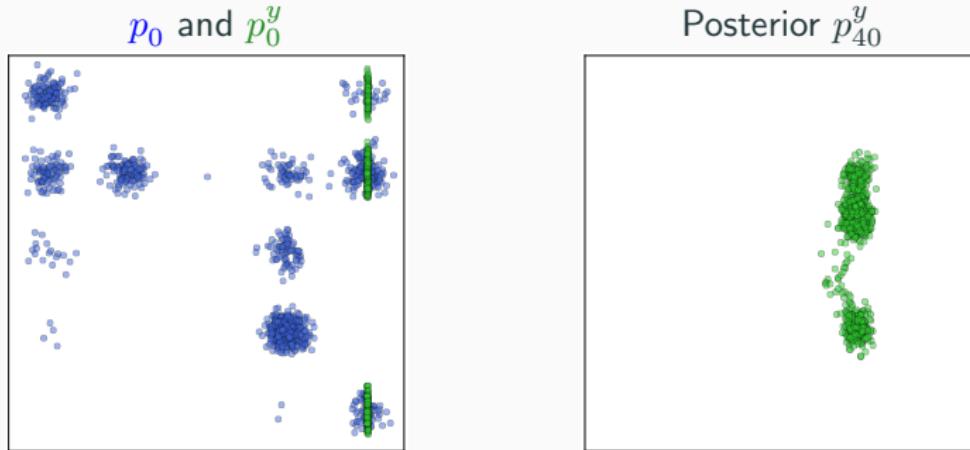
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 11:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

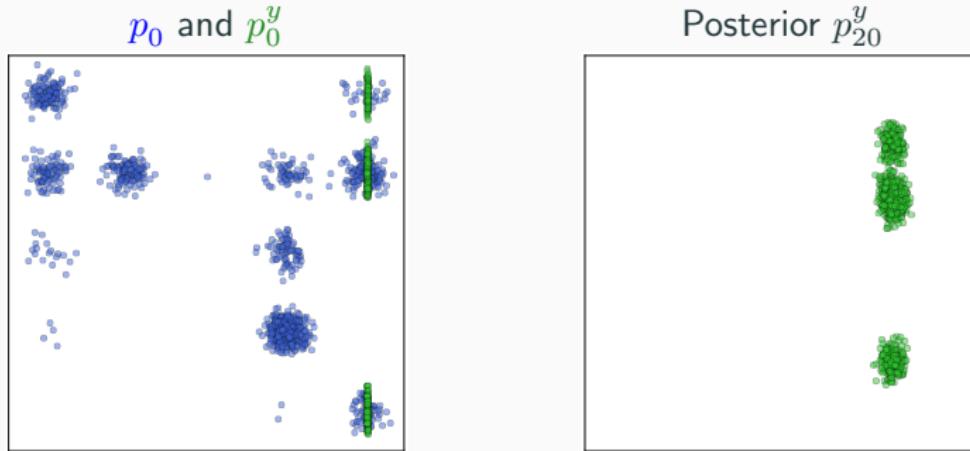
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 12:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

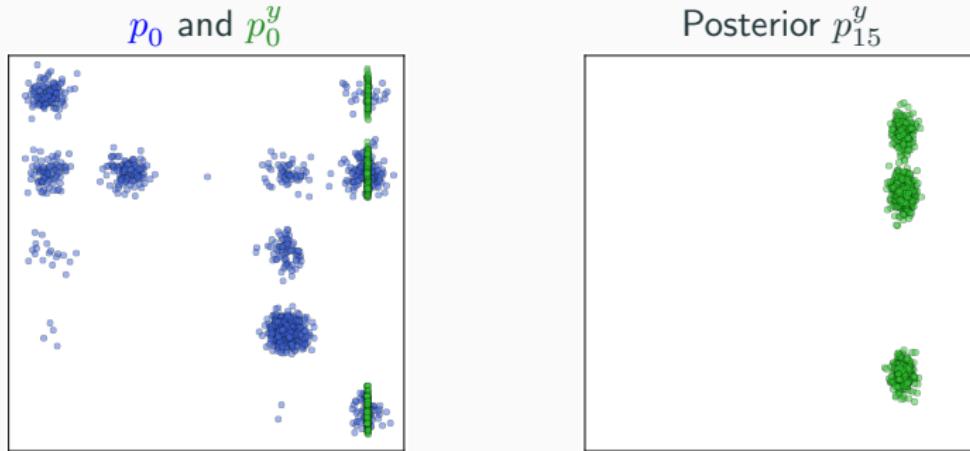
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 13:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

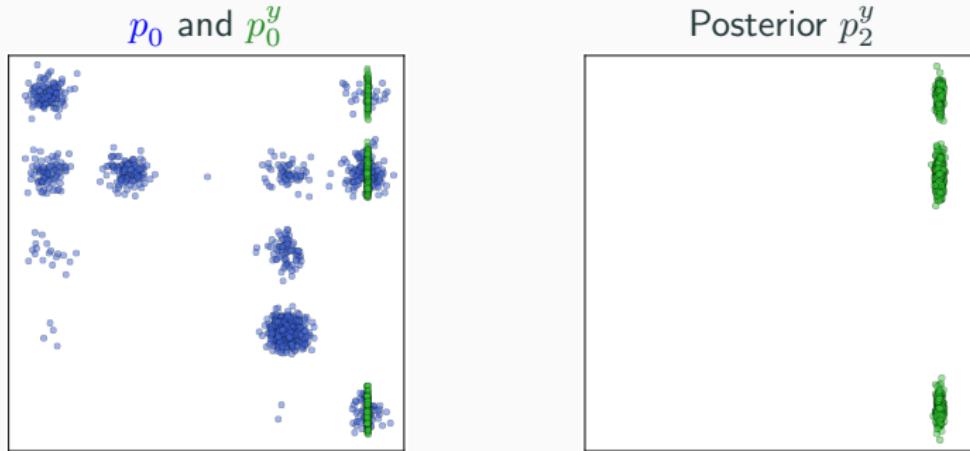
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 14:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

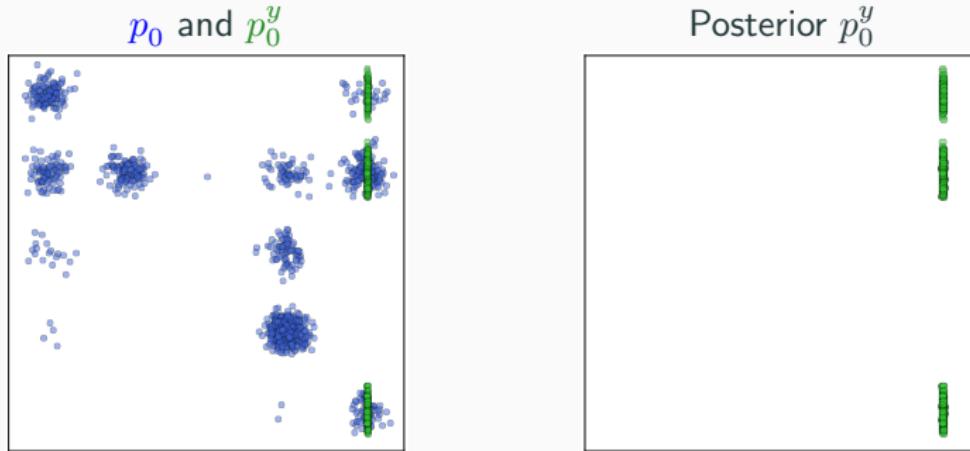
$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 15:** **Left plot:** samples from the prior  $p_0$  and posterior  $p_0^y$ . **Right plot:** samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Illustration

$\rightsquigarrow \{p_k^y\}_{k=1}^n$  is available in closed form for the Gaussian mixture example.



**Figure 16:** Left plot: samples from the prior  $p_0$  and posterior  $p_0^y$ . Right plot: samples from the posterior proposals  $p_k^y$  for time steps ranging from  $n := 500$  to 0.

# Toy examples

- ~ 25 Gaussian mixture example with means

$$\mu_{i,j} = (8i, 8j, \dots, 8i, 8j), \quad (i, j) \in \{-2, \dots, 2\}$$

with unit covariance matrices. We randomly draw the weights of the mixture and the forward operator  $A$  and  $\sigma_y$  for the inverse problem  $\leadsto \nabla \log p_k$  is available in **closed form**.

- ~ 20 component mixture of translated and rotated Funnel distributions. We learn the score and consider the ground truth to be samples from parallel NUTS with very long chains.

# Toy examples

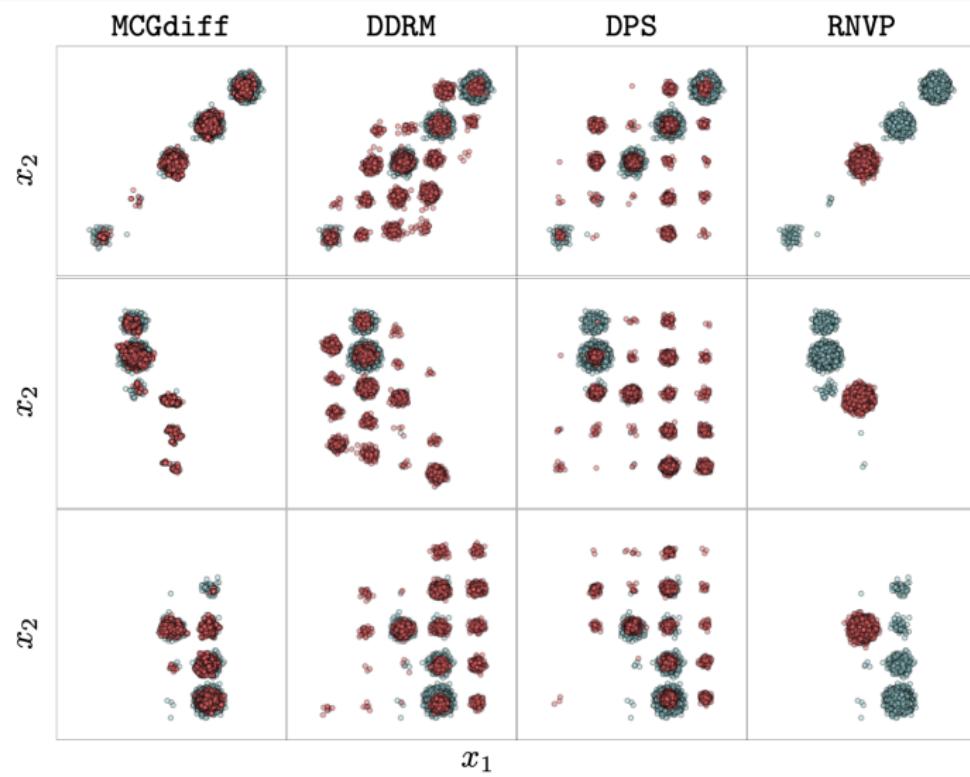
| $d$ | $d_y$ | MCGdiff                           | DDRM            | DPS             | RNVP            |
|-----|-------|-----------------------------------|-----------------|-----------------|-----------------|
| 80  | 1     | <b><math>1.39 \pm 0.45</math></b> | $5.64 \pm 1.10$ | $4.98 \pm 1.14$ | $6.86 \pm 0.88$ |
| 80  | 2     | <b><math>0.67 \pm 0.24</math></b> | $7.07 \pm 1.35$ | $5.10 \pm 1.23$ | $7.79 \pm 1.50$ |
| 80  | 4     | <b><math>0.28 \pm 0.14</math></b> | $7.81 \pm 1.48$ | $4.28 \pm 1.26$ | $7.95 \pm 1.61$ |
| 800 | 1     | <b><math>2.40 \pm 1.00</math></b> | $7.44 \pm 1.15$ | $6.49 \pm 1.16$ | $7.74 \pm 1.34$ |
| 800 | 2     | <b><math>1.31 \pm 0.60</math></b> | $8.95 \pm 1.12$ | $6.88 \pm 1.01$ | $8.75 \pm 1.02$ |
| 800 | 4     | <b><math>0.47 \pm 0.19</math></b> | $8.39 \pm 1.48$ | $5.51 \pm 1.18$ | $7.81 \pm 1.63$ |

| $d$ | $d_y$ | MCGdiff                           | DDRM            | DPS             | RNVP            |
|-----|-------|-----------------------------------|-----------------|-----------------|-----------------|
| 6   | 1     | <b><math>1.95 \pm 0.43</math></b> | $4.20 \pm 0.78$ | $5.43 \pm 1.05$ | $6.16 \pm 0.65$ |
| 6   | 3     | <b><math>0.73 \pm 0.33</math></b> | $2.20 \pm 0.67$ | $3.47 \pm 0.78$ | $4.70 \pm 0.90$ |
| 6   | 5     | <b><math>0.41 \pm 0.12</math></b> | $0.91 \pm 0.43$ | $2.07 \pm 0.63$ | $3.52 \pm 0.93$ |
| 10  | 1     | <b><math>2.45 \pm 0.42</math></b> | $3.82 \pm 0.64$ | $4.30 \pm 0.91$ | $6.04 \pm 0.38$ |
| 10  | 3     | <b><math>1.07 \pm 0.26</math></b> | $4.94 \pm 0.87$ | $5.38 \pm 0.84$ | $5.91 \pm 0.64$ |
| 10  | 5     | <b><math>0.71 \pm 0.12</math></b> | $2.32 \pm 0.74$ | $3.74 \pm 0.77$ | $5.11 \pm 0.69$ |

**Figure 17:** Sliced Wasserstein between samples of the target posterior and the empirical measure returned by each method. **Top:** Gaussian mixture. **Bottom:** Funnel mixture. We show the 95% CLT interval over 20 seeds.

# Toy examples



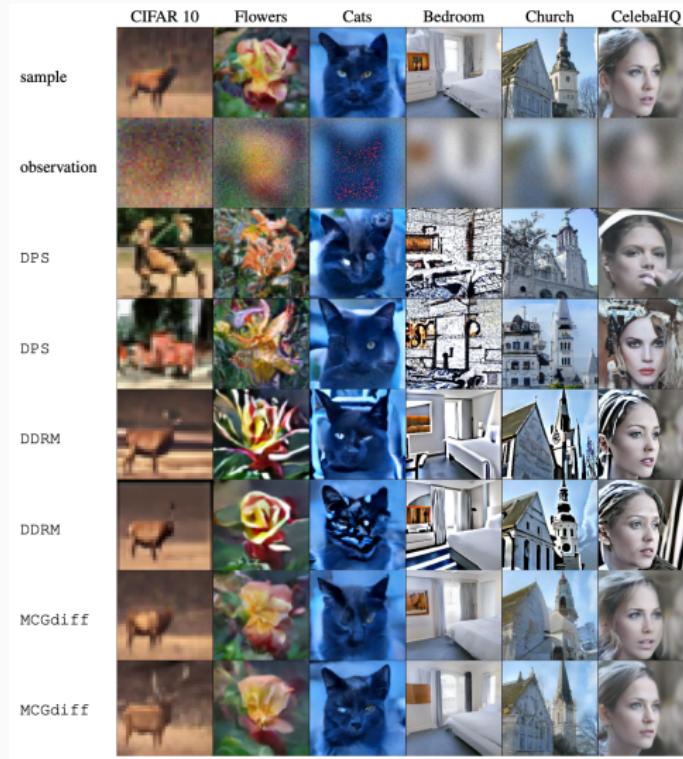
# Imaging experiments

- ~~ Diffusion models learned on different datasets of image sizes varying from  $(64, 64, 3)$  to  $(256, 256, 3)$ .
- ~~ We run parallel SMCs with **N = 64** particles.

# Super-resolution example



# Deblurring example



# Inpainting example



## Divide-and-conquer posterior sampling

---

# Sequence of distributions

Let  $(k_\ell)_{\ell=0}^L$  be an increasing sequence in  $\llbracket 0, n \rrbracket$  with  $k_0 = 0$  and  $k_L = n$ .

Consider

$$p_{k_\ell}^y(\mathrm{d}x_{k_\ell}) \propto g_{k_\ell}^y(x_{k_\ell}) p_{k_\ell}(\mathrm{d}x_\ell),$$

with

$$g_{k_\ell}^y(x_{k_\ell}) = \mathcal{N}(\sqrt{\alpha_{k_\ell}} y; Ax_{k_\ell}, \sigma_{y,k_\ell}^2 \mathbf{I}_{d_y}).$$

- $L$  is typically much smaller than  $n$ .
- This is the same sequence of distribution as in our SMC approach but now we only consider a **small number  $L$**  of intermediate distributions.
- Our goal is to recursively sample from each one of them without having to evolve  **$N$  particles** in parallel.
- We also want to solve the “image inconsistency” problem observed in our **SMC method**.

# Recursion

Since

$$p_{k_\ell}(\mathrm{d}x_{k_\ell}) = \int \left\{ \prod_{j=k_\ell}^{k_{\ell+1}-1} p_{j|j+1}(\mathrm{d}x_j|x_{j+1}) \right\} p_{k_{\ell+1}}(\mathrm{d}x_{k_{\ell+1}}),$$

we can write  $p_{k_\ell}^y$  in terms of forward smoothing kernels, i.e.

$$p_{k_\ell}^y(\mathrm{d}x_{k_\ell}) = \int \left\{ \prod_{j=k_\ell}^{k_{\ell+1}-1} p_{j|j+1}^{y,\ell}(\mathrm{d}x_j|x_{j+1}) \right\} p_{k_{\ell+1}}^{y,\ell}(\mathrm{d}x_{k_{\ell+1}})$$

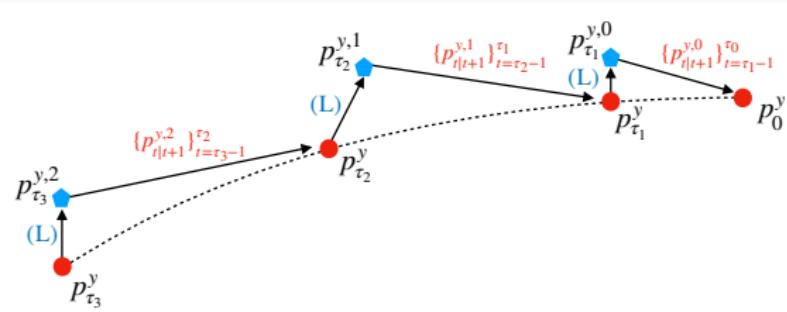
where

$$\begin{aligned} p_{k_{\ell+1}}^{y,\ell}(\mathrm{d}x_{k_{\ell+1}}) &\propto \beta_{k_\ell|k_{\ell+1}}^{y,\ell}(x_{k_{\ell+1}}) p_{k_{\ell+1}}(\mathrm{d}x_{k_{\ell+1}}), \\ p_{j|j+1}^{y,\ell}(\mathrm{d}x_j|x_{j+1}) &\propto \beta_{k_\ell|j}^{y,\ell}(x_j) p_{j|j+1}(\mathrm{d}x_j|x_{j+1}), \end{aligned}$$

and for all  $j \in \llbracket k_\ell, k_{\ell+1} \rrbracket$

$$\beta_{k_\ell|j}^{y,\ell}(x_j) := \int g_{k_\ell}^y(x_{k_\ell}) p_{k_\ell|j}(\mathrm{d}x_{k_\ell}|x_j).$$

# DCPS summary



**Figure 18:** Illustration of idealized DCPS.

Starting at an approximate sample  $X_{k_{\ell+1}}^y$  from  $p_{k_{\ell+1}}^y$

- Use ULA initialized at  $X_{k_{\ell+1}}^y$  to obtain an approximate sample from  $X_{k_{\ell+1}}^{y,\ell}$ .
- Starting from  $X_{k_{\ell+1}}^{y,\ell}$ , simulate a Markov chain with transition kernels  $(p_{j|j+1}^{y,\ell})_{j=k_{\ell+1}-1}^{k_{\ell}}$
- Repeat until the posterior of interest is reached.

# Backward function approximation

- The first source of intractability are the backward functions  $\beta_{k_\ell|j}^{y,\ell}$ .
- This is the same problem as before, however note that now they are expressed as an integral under  $p_{k_\ell|j}(\cdot|x_j)$  with  $j \in [k_\ell + 1, k_{\ell+1}]$  instead of  $p_{0|j}(\cdot|x_j)$  for  $j \in [0, n]$ .
- This is more convenient since we expect Gaussian approximations of  $p_{k_\ell|j}(\cdot|x_j)$  to be more accurate than those of  $p_{0|j}(\cdot|x_j)$ .

# Backward kernel approximation

Assume again that **forward=backward**. Then for  $j \in [k_\ell + 1, k_{\ell+1}]$ ,

$$p_{k_\ell|j}(\mathrm{d}x_{k_\ell}|x_j) = \int q_{k_\ell|0,j}(\mathrm{d}x_{k_\ell}|x_0, x_j) p_{0|j}(\mathrm{d}x_0|x_j),$$

Let  $\hat{p}_{0|j}(\cdot|x_j)$  be an approximation of  $p_{0|j}(\cdot|x_j)$  and define

$$\hat{p}_{k_\ell|j}(\mathrm{d}x_{k_\ell}|x_j) = \int q_{k_\ell|0,j}(\mathrm{d}x_{k_\ell}|x_0, x_j) \hat{p}_{0|j}(\mathrm{d}x_0|x_j)$$

- For DPS (Chung et al., 2023),  $\hat{p}_{0|j}(\mathrm{d}x_0|x_j) = \delta_{\hat{x}_{0|j}^\theta(x_j)}(\mathrm{d}x_0)$ .
- For Song et al. (2023),  $\hat{p}_{0|j}(\mathrm{d}x_0|x_j) = \mathcal{N}(\mathrm{d}x_0; \hat{x}_{0|j}^\theta(x_j), r_j^2 \mathbf{I}_{d_y})$ .
- In both cases,  $\hat{p}_{k_\ell|j}(\cdot|x_j)$  is computable in **closed form**. We write

$$\hat{p}_{k_\ell|j}(\mathrm{d}x_{k_\ell}|x_j) = \mathcal{N}(\mathrm{d}x_{k_\ell}; \mu_{k_\ell|j}(x_j), \sigma_{k_\ell|j}^2 \mathbf{I}_{d_x}).$$

where both the mean and variance depend on the approximation used.

# Backward kernel approximation

## Proposition

Assume **forward=backward**. For all  $\ell \in \llbracket 0, L \rrbracket$ ,  $j \in \llbracket k_\ell + 1, k_{\ell+1} \rrbracket$ ,

$$W_2(\hat{p}_{k_\ell|j}(\cdot|x_j), p_{k_\ell|j}(\cdot|x_j)) \leq \frac{\sqrt{\alpha_{k_\ell}}(1 - \alpha_j/\alpha_{k_\ell})}{1 - \alpha_j} W_2(\hat{p}_{0|j}(\cdot|x_j), p_{0|j}(\cdot|x_j)).$$

where  $\frac{\sqrt{\alpha_{k_\ell}}(1 - \alpha_j/\alpha_{k_\ell})}{1 - \alpha_j} < 1$  and goes to 0 as  $j \rightarrow k_\ell$ .

- We improve upon the previous approximations by performing Gaussian approximations on intervals  $\llbracket k_\ell, k_{\ell+1} \rrbracket$  of moderate size.
- Our approximation of the backward function is then

$$\begin{aligned}\beta_{k_\ell|j}^{y,\ell}(x_j) &\approx \hat{\beta}_{k_\ell|j}^{y,\ell}(x_j) := \int g_{k_\ell}^y(x_{k_\ell}) \hat{p}_{k_\ell|j}(\mathrm{d}x_{k_\ell}|x_j) \\ &= \mathcal{N}(\sqrt{\alpha_{k_\ell}} y; A\mu_{k_\ell|j}(x_j), \sigma_{k_\ell|j}^2 AA^\top + \sigma_{y,\ell}^2 I_{d_y}).\end{aligned}$$

# FSK approximation

Recall that the quantities of interest are

$$\begin{aligned} p_{j|j+1}^{y,\ell}(\mathrm{d}x_j|x_{j+1}) &\propto \beta_{k_\ell|j}^{y,\ell}(x_j) p_{j|j+1}(\mathrm{d}x_j|x_{j+1}), \\ p_{k_{\ell+1}}^{y,\ell}(\mathrm{d}x_{k_{\ell+1}}) &\propto \beta_{k_\ell|k_{\ell+1}}^{y,\ell}(x_{k_{\ell+1}}) p_{k_{\ell+1}}(\mathrm{d}x_{k_{\ell+1}}). \end{aligned}$$

Given the previous approximation of the backward function, we replace them instead with

$$\begin{aligned} \hat{p}_{j|j+1}^{y,\ell}(\mathrm{d}x_j|x_{j+1}) &\propto \hat{\beta}_{k_\ell|j}^{y,\ell}(x_j) p_{j|j+1}(\mathrm{d}x_j|x_{j+1}), \\ \hat{p}_{k_{\ell+1}}^{y,\ell}(\mathrm{d}x_{k_{\ell+1}}) &\propto \hat{\beta}_{k_\ell|k_{\ell+1}}^{y,\ell}(x_{k_{\ell+1}}) p_{k_{\ell+1}}(\mathrm{d}x_{k_{\ell+1}}), \end{aligned}$$

- Still, while now we can evaluate the density  $\hat{p}_{j|j+1}^{y,\ell}(\cdot|x_{j+1})$  we still **cannot sample** from it.
- We can approximately sample from  $\hat{p}_{k_{\ell+1}}^{y,\ell}$  using ULA.

# Variational approximation I

For a **fixed**  $x_{j+1}$  we seek a **mean-field Gaussian variational approximation** of  $\hat{p}_{j|j+1}^{y,\ell}(\cdot|x_{j+1})$  by solving

$$\operatorname{argmin}_{r_{j|j+1}^{y,\ell}(\cdot|x_{j+1}) \in \mathcal{G}_D} \text{KL}(r_{j|j+1}^{y,\ell}(\cdot|x_{j+1}) \parallel \hat{p}_{j|j+1}^{y,\ell}(\cdot|x_{j+1})),$$

where  $\mathcal{G}_D := \{\mathcal{N}(\mu, \operatorname{diag}(\sigma)) : \mu \in \mathbb{R}^{d_x}, \sigma \in \mathbb{R}_{>0}^{d_x}\}$ .

- We only learn vectors  $(\mu, \sigma)$  that depend on the value of  $X_{j+1}^{y,\ell}$  and do not seek to generalize as this incurs **problem dependent, heavy training**.

# Variational approximation II

Letting  $r_{j|j+1}^{y,\ell}(\cdot|X_{j+1}^{y,\ell}) = \mathcal{N}(\mu_{j|j+1}^{y,\ell}, \text{diag}(\text{e}^{s_{j|j+1}^{y,\ell}}))$  where  $s_{j|j+1}^{y,\ell} \in \mathbb{R}^{d_x}$ ,

$$\begin{aligned} & \text{KL}(r_{j|j+1}^{y,\ell}(\cdot|X_{j+1}^{y,\ell}) \parallel \hat{p}_{j|j+1}^{y,\ell}(\cdot|X_{j+1}^{y,\ell})) \\ &= -\mathbb{E}[\log \hat{\beta}_{k_\ell|j}^{y,\ell}(\mu_{j|j+1}^{y,\ell} + \text{diag}(\text{e}^{s_{j|j+1}^{y,\ell}})Z)] + \frac{\|\mu_{j|j+1}^{y,\ell} - \mu_{j|j+1}(X_{j+1}^{y,\ell})\|^2}{2\sigma_{m|m+1}^2} \\ & \quad - \frac{1}{2} \sum_{i=1}^{d_x} \left( s_{j|j+1,i}^{y,\ell} - \frac{\text{e}^{s_{j|j+1,i}^{y,\ell}}}{\sigma_{m|m+1}^2} \right), \end{aligned}$$

- We perform the optimization using SGD.
- Crucially, we normalize the gradients to ensure the stability of the training procedure.
- In practice, we only perform **2 or 3** SGD steps.

# Tamed ULA steps

We now turn to the Langevin steps on  $\hat{p}_{k_{\ell+1}}^{y,\ell}$ .

As the marginals  $(p_k)_{k=0}^n$  approximate the true marginals of the forward process initialized at the data distribution  $\pi$ , we may use

$$s_k^\theta(x_k) = -(x_k - \sqrt{\alpha_k} \hat{x}_{0|k}^\theta(x_k)) / (1 - \alpha_k),$$

as a substitute for  $\nabla_{x_k} \log p_k(x_k)$ , following [Dhariwal and Nichol \(2021\)](#).

We sample approximately from  $\hat{p}_{k_{\ell+1}}^{y,\ell}$  by running  $M$  steps of the Tamed Unadjusted Langevin scheme ([Brosse et al., 2019](#))

$$X_{j+1} = X_j + \gamma G_\gamma^{y,\ell}(X_j) + \sqrt{2\gamma} Z_j, \quad X_0 = X_{k_{\ell+1}}^y, \quad (1)$$

where

$$G_\gamma^{y,\ell}(x) := \frac{\nabla \log \hat{\beta}_{k_\ell|k_{\ell+1}}^{y,\ell}(x) + s_{k_{\ell+1}}^\theta(x)}{1 + \gamma \|\nabla \log \hat{\beta}_{k_\ell|k_{\ell+1}}^{y,\ell}(x) + s_{k_{\ell+1}}^\theta(x)\|},$$

and set  $X_{k_{\ell+1}}^{y,\ell} := X_M$ .

# Summary

Given an approximate sample  $X_{k_{\ell+1}}^y$  from  $\hat{p}_{k_{\ell+1}}^y$ ,

- Run TULA starting from  $X_{k_{\ell+1}}^y$  to obtain  $X_{k_{\ell+1}}^{y,\ell}$  approximately distributed according  $\hat{p}_{k_{\ell+1}}^{y,\ell}$ .
- Sample  $(X_j^{y,\ell})_{j=k_{\ell+1}}^{k_{\ell}}$ : given  $X_{j+1}^{y,\ell}$  with  $j \in [k_{\ell}, k_{\ell+1} - 1]$ ,
  - Find variational approximation  $r_{j|j+1}^{y,\ell}(\cdot | X_{j+1}^{y,\ell})$ .
  - Draw  $X_j^{y,\ell} \sim r_{j|j+1}^{y,\ell}(\cdot | X_{j+1}^{y,\ell})$ .
- Repeat these steps.

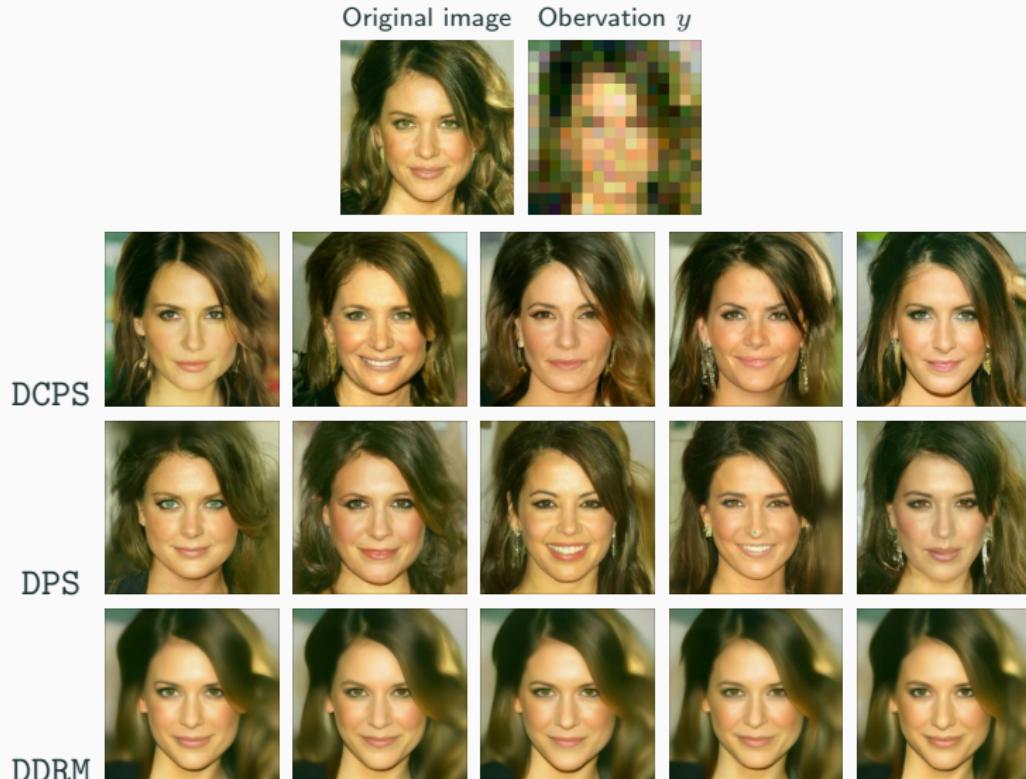
# Toy experiments

- Same 25 Gaussian mixture example.
- $M$  Langevin steps at the beginning of each block and  $L = 4$ .

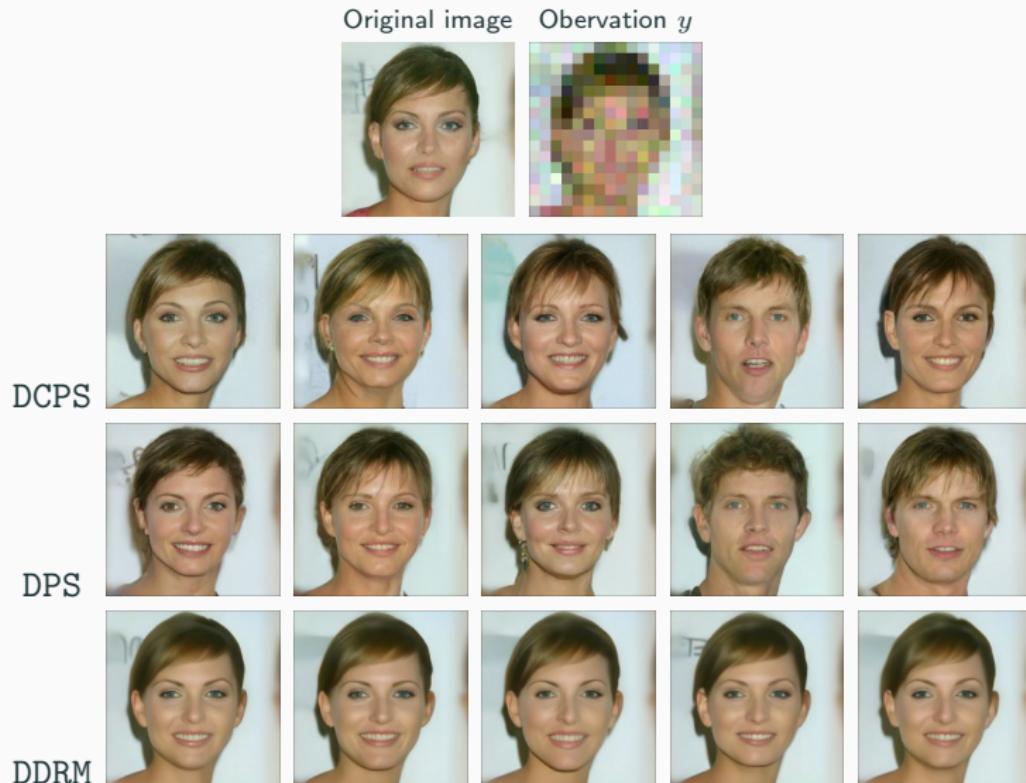
| $d_x = 10, d_y = 1$ |                      |                 | $d_x = 100, d_y = 1$ |                      |                 |
|---------------------|----------------------|-----------------|----------------------|----------------------|-----------------|
|                     | SW                   | $\Delta w$      |                      | SW                   | $\Delta w$      |
| MCGDiff             | 2.25/2.69 $\pm$ 2.07 | 0.32 $\pm$ 0.20 |                      | 2.72/3.13 $\pm$ 1.76 | 0.42 $\pm$ 0.19 |
| DPS                 | 3.12/5.64 $\pm$ 8.45 | 0.20 $\pm$ 0.12 |                      | 4.29/4.93 $\pm$ 4.85 | 0.35 $\pm$ 0.25 |
| DDRM                | 2.66/3.06 $\pm$ 1.90 | 0.36 $\pm$ 0.16 |                      | 5.97/6.26 $\pm$ 2.33 | 0.52 $\pm$ 0.19 |
| DCPS <sub>50</sub>  | 1.95/2.70 $\pm$ 2.28 | 0.17 $\pm$ 0.25 |                      | 4.40/4.72 $\pm$ 2.18 | 0.44 $\pm$ 0.16 |
| DCPS <sub>500</sub> | 1.26/2.59 $\pm$ 2.83 | 0.13 $\pm$ 0.30 |                      | 2.81/3.22 $\pm$ 2.21 | 0.32 $\pm$ 0.18 |

**Table 1:** Results for the Gaussian mixture experiment. Results for the SW distance are shown in median/mean  $\pm$  standard deviation format.

# Super-resolution experiments



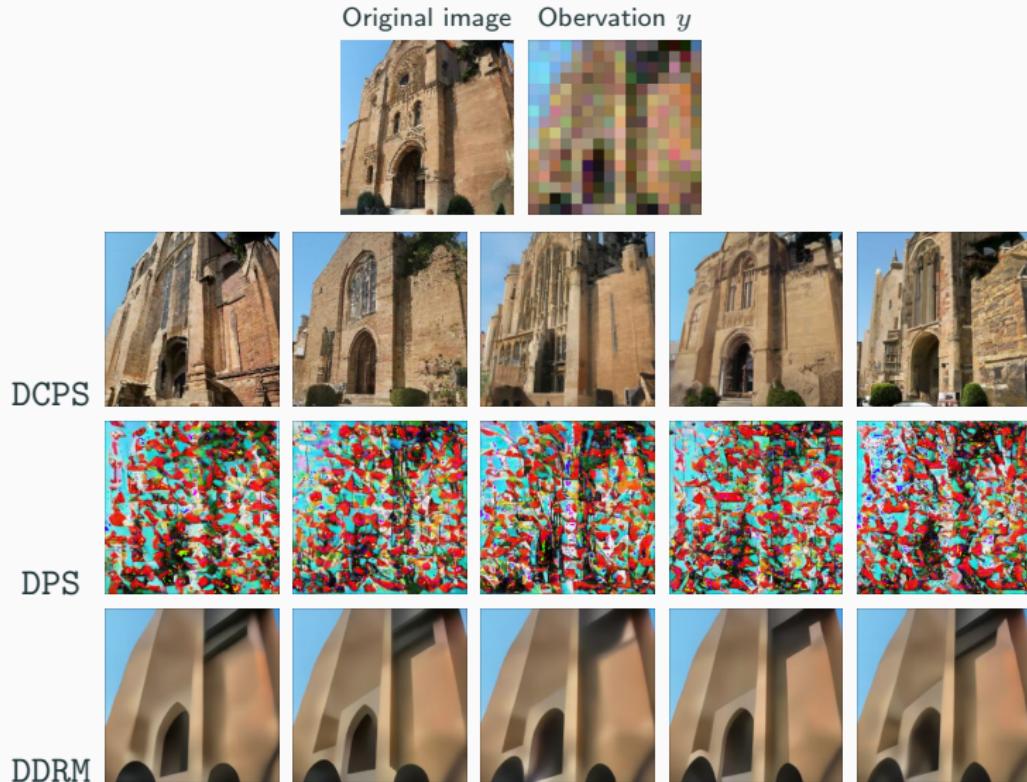
# Super-resolution experiments



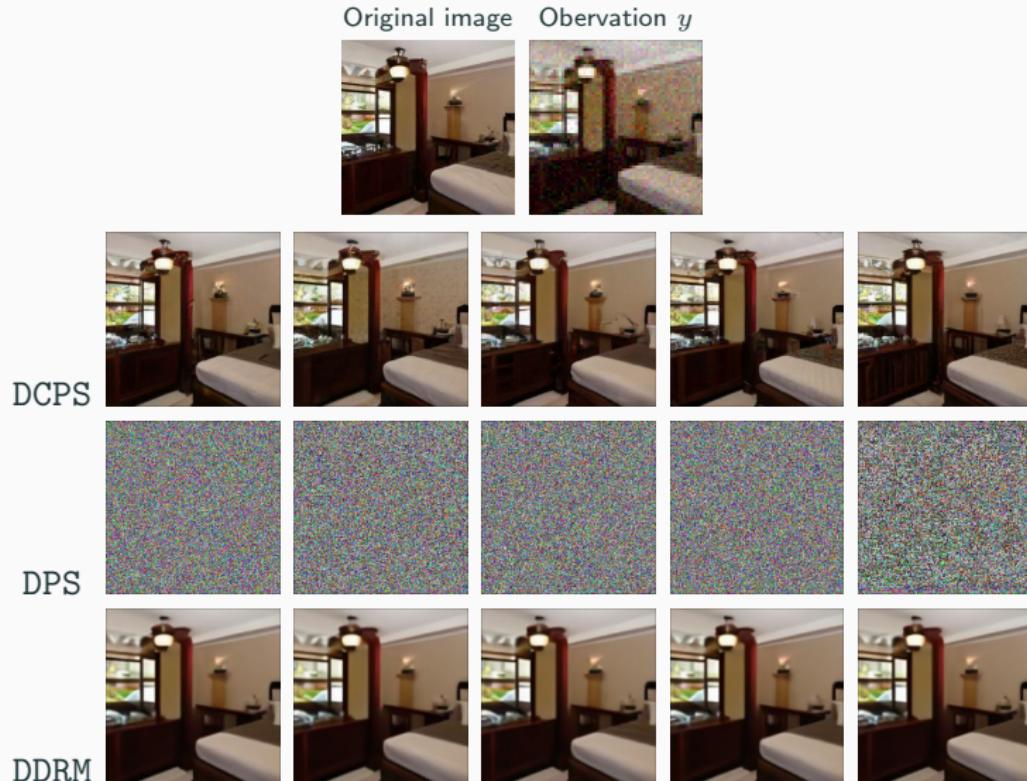
# Super-resolution experiments



# Super-resolution experiments



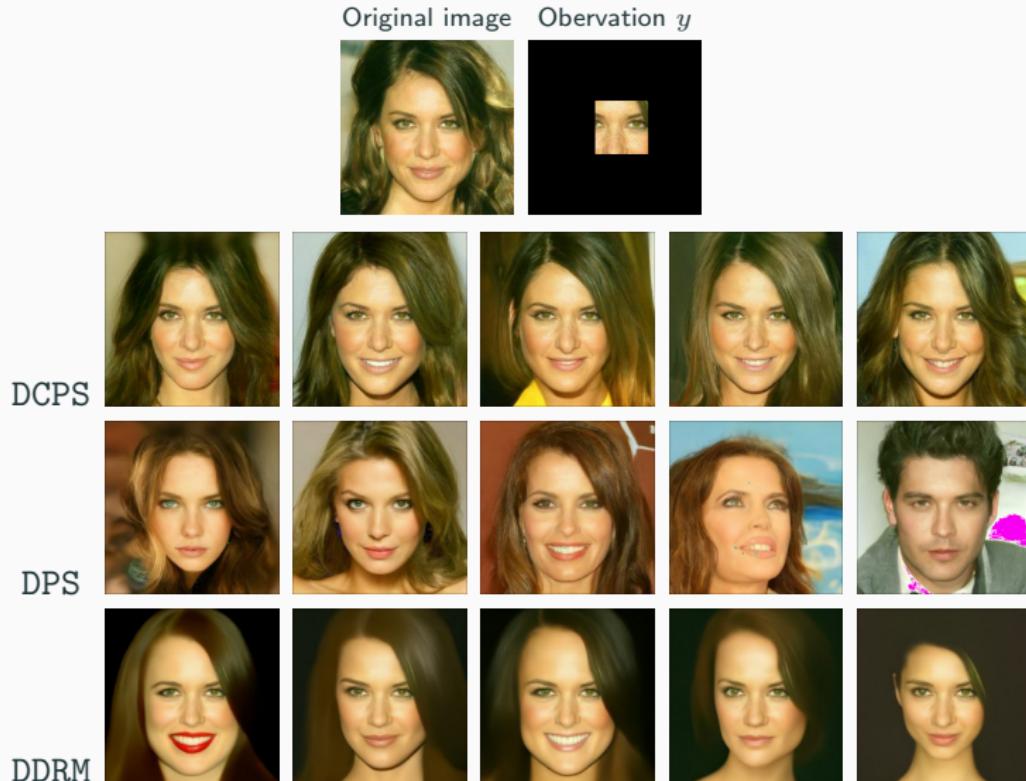
# Super-resolution experiments



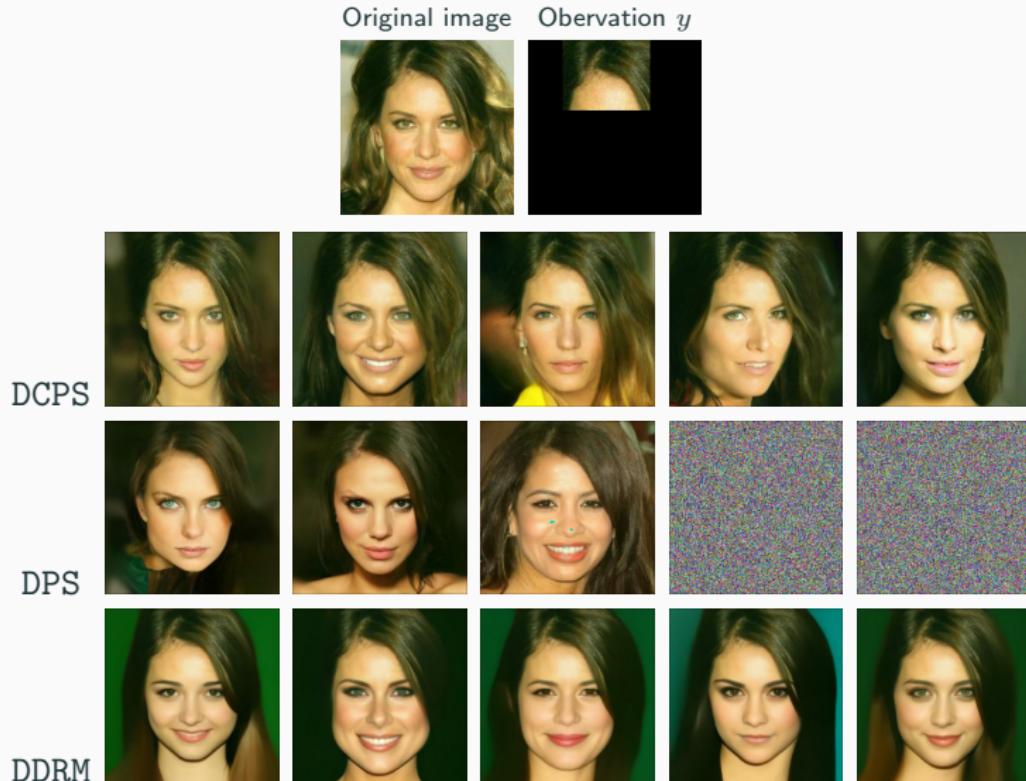
# Super-resolution experiments



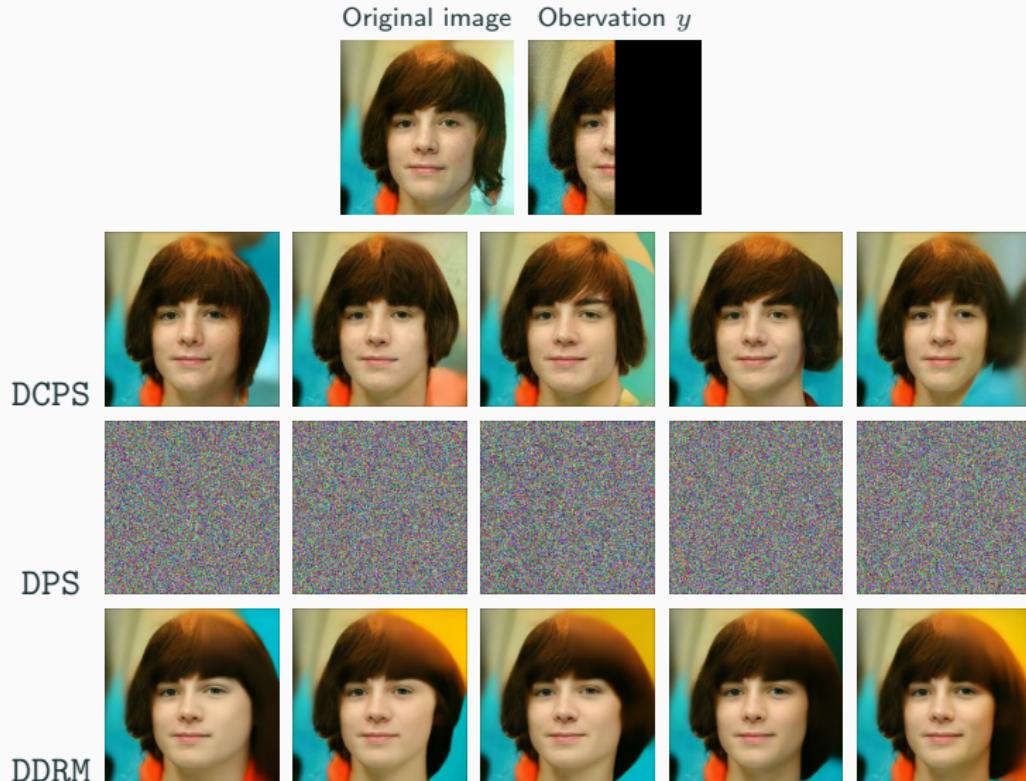
# Inpainting and outpainting experiments



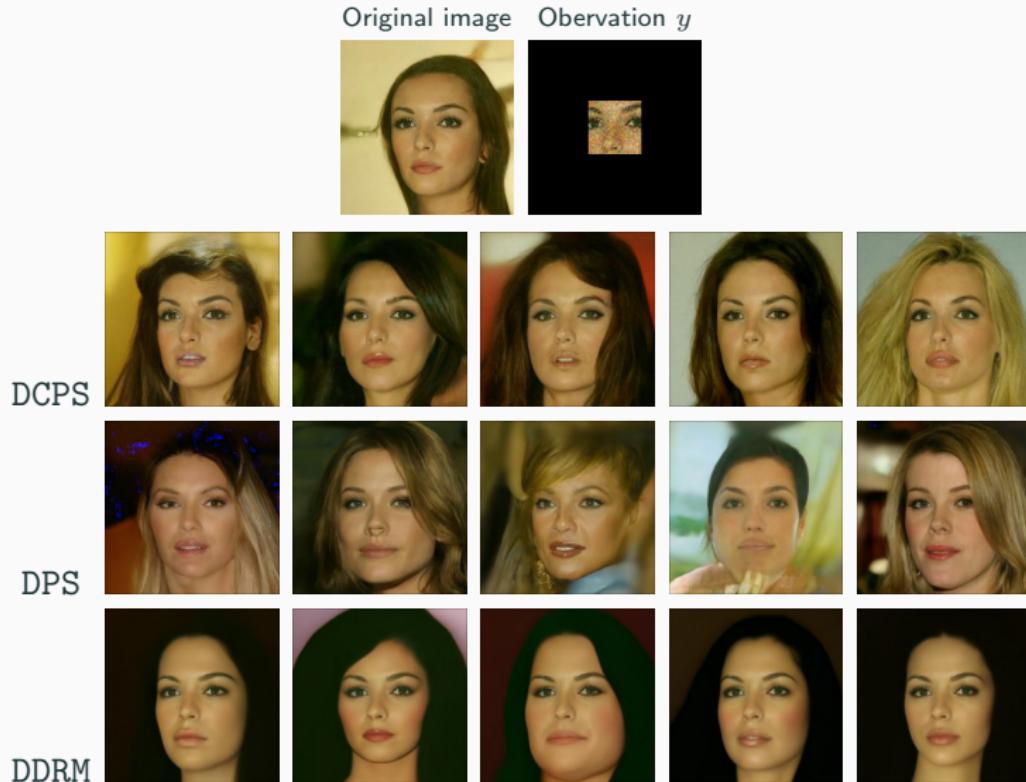
# Inpainting and outpainting experiments



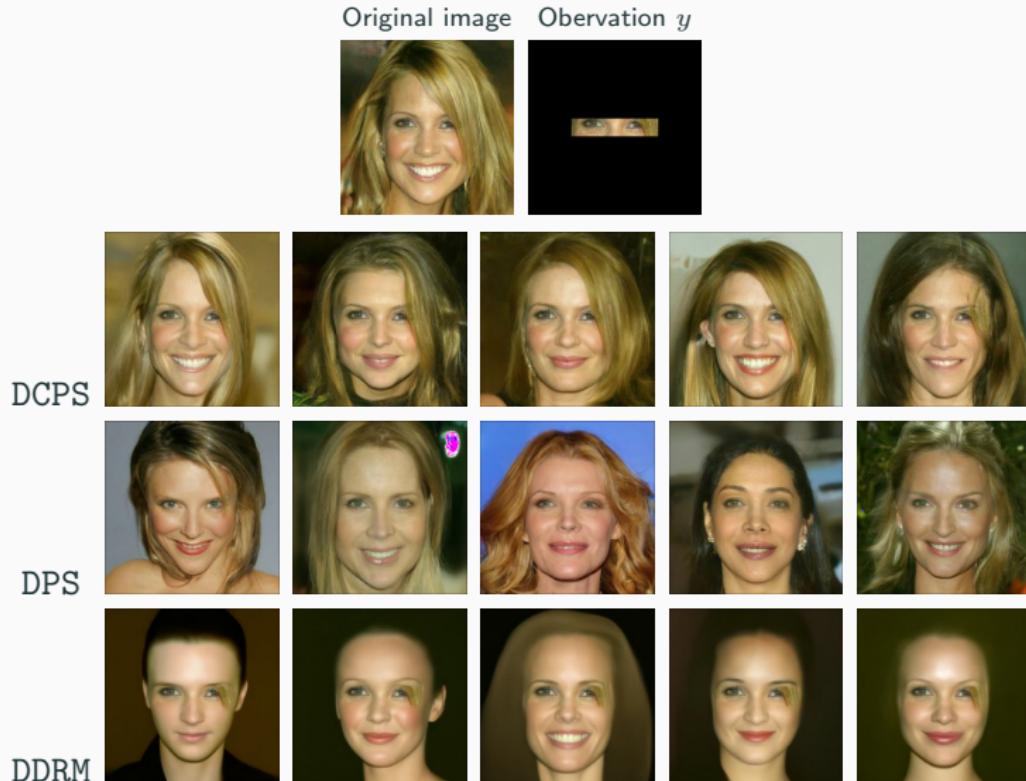
# Inpainting and outpainting experiments



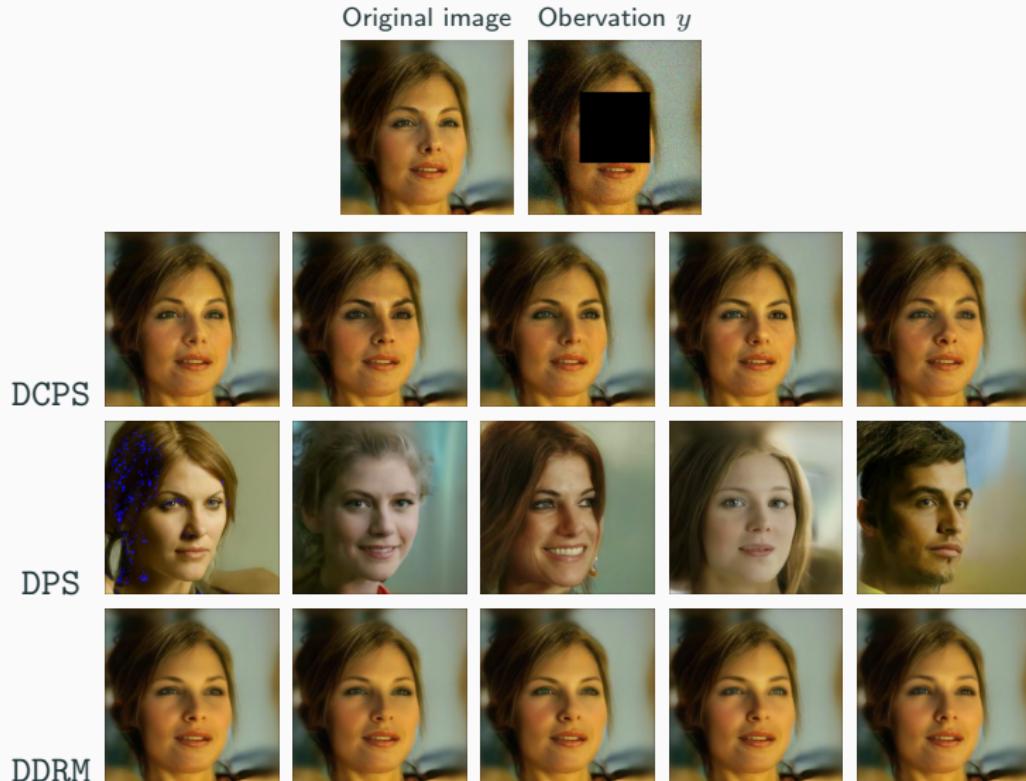
# Inpainting and outpainting experiments



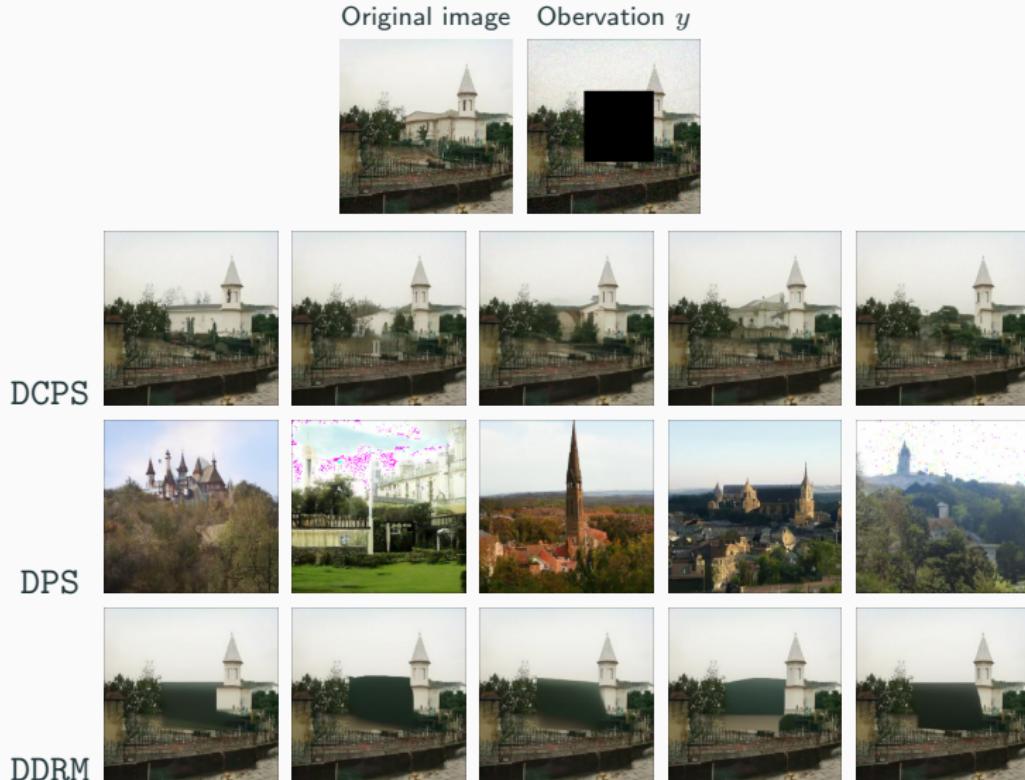
# Inpainting and outpainting experiments



# Inpainting and outpainting experiments



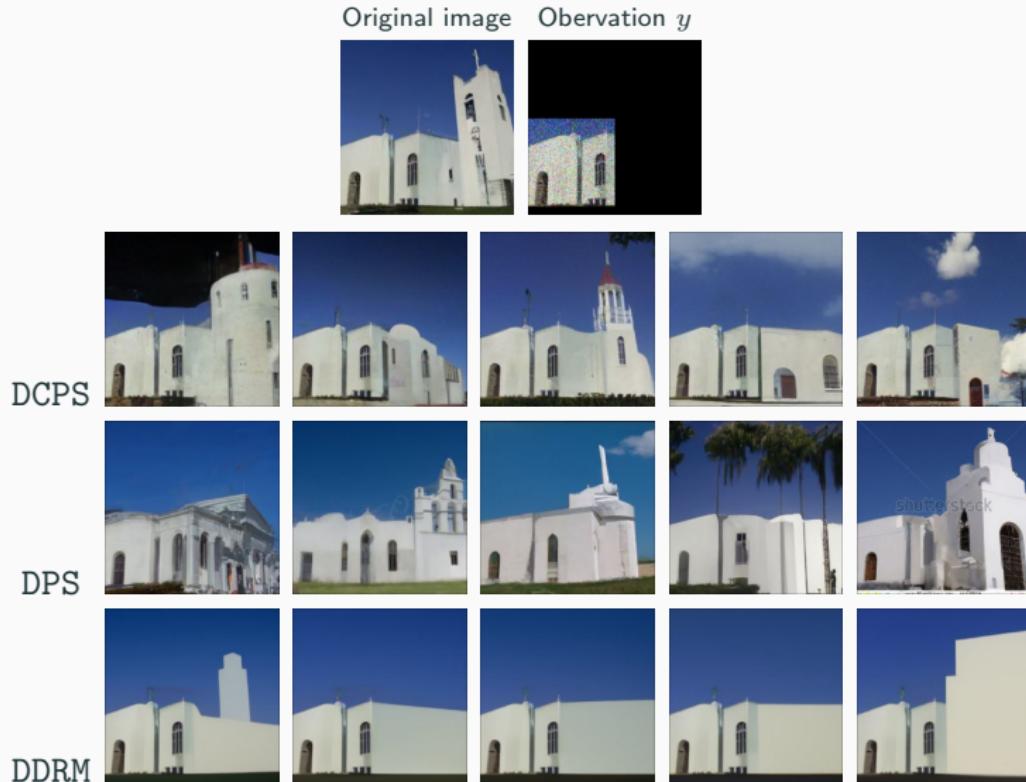
# Inpainting and outpainting experiments



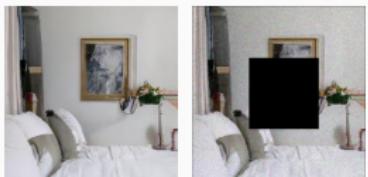
# Inpainting and outpainting experiments



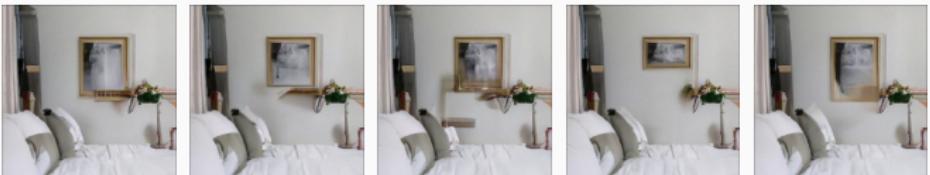
# Inpainting and outpainting experiments



# Inpainting and outpainting experiments

Original image   Observation  $y$ 

DCPS



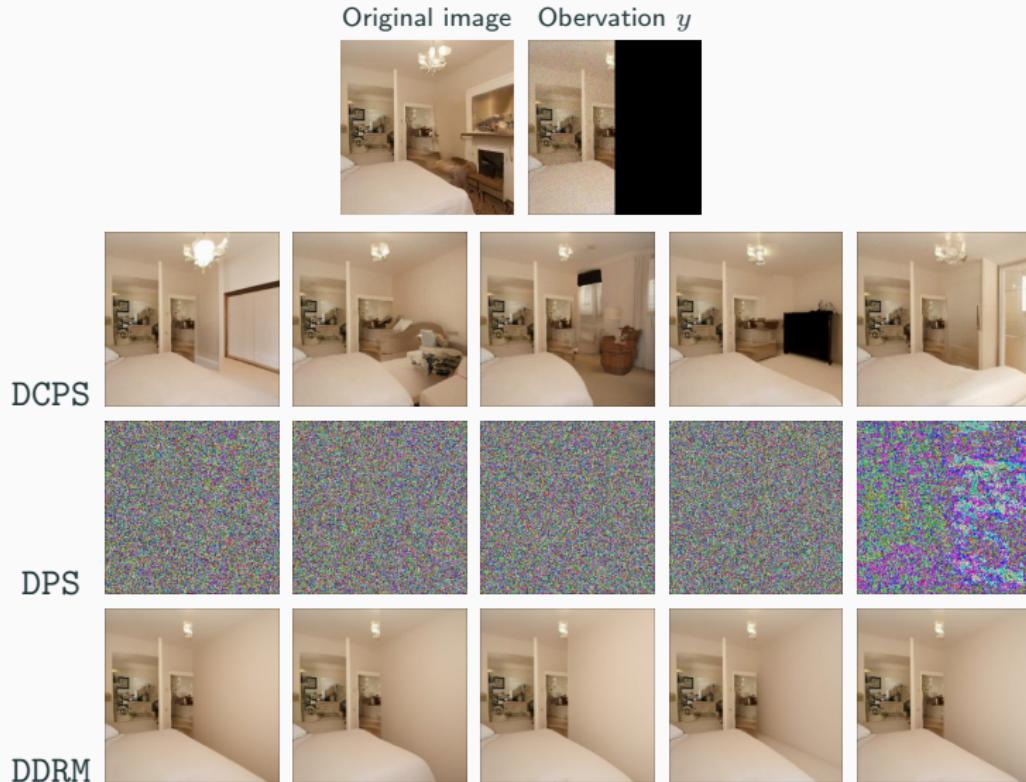
DPS



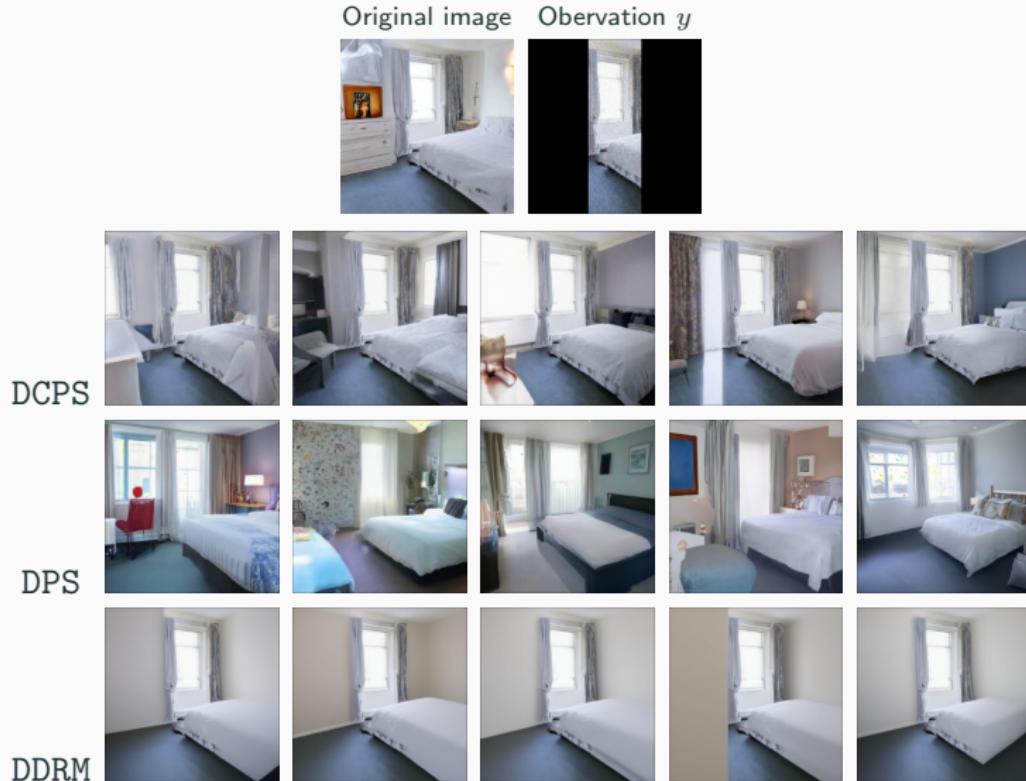
DDRM



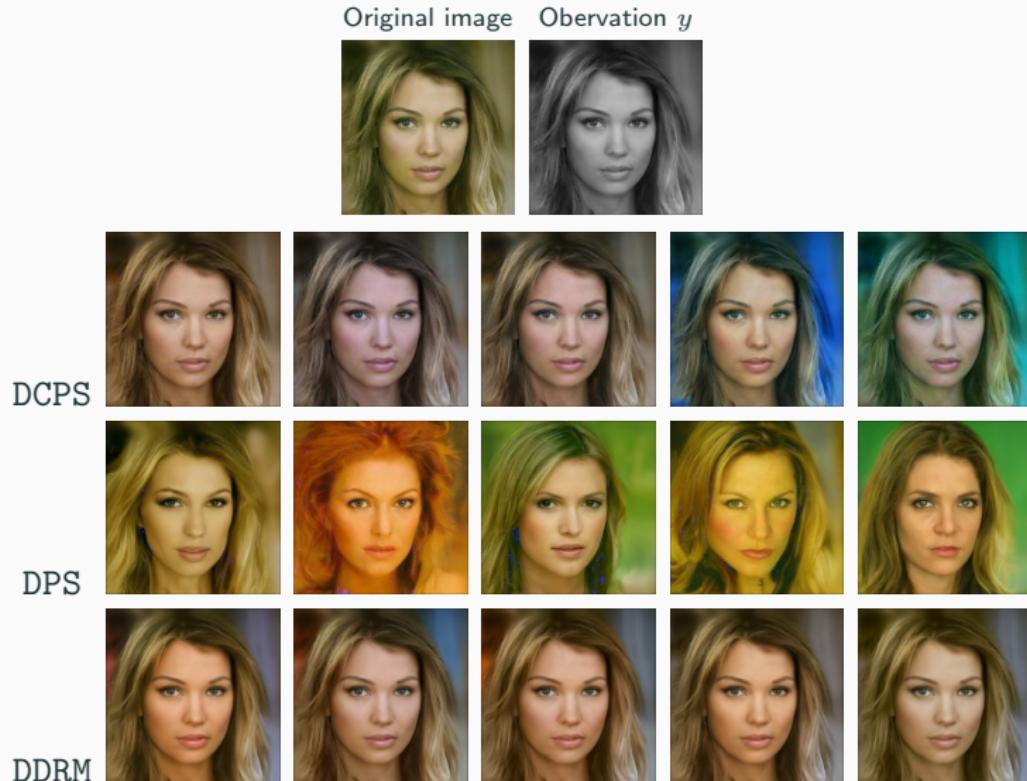
# Inpainting and outpainting experiments



# Inpainting and outpainting experiments



# Colorization experiments



# Colorization experiments

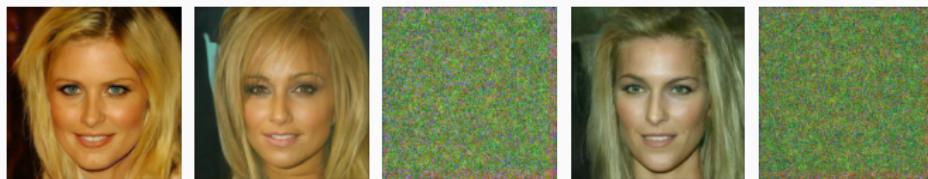
Original image   Observation  $y$



DCPS



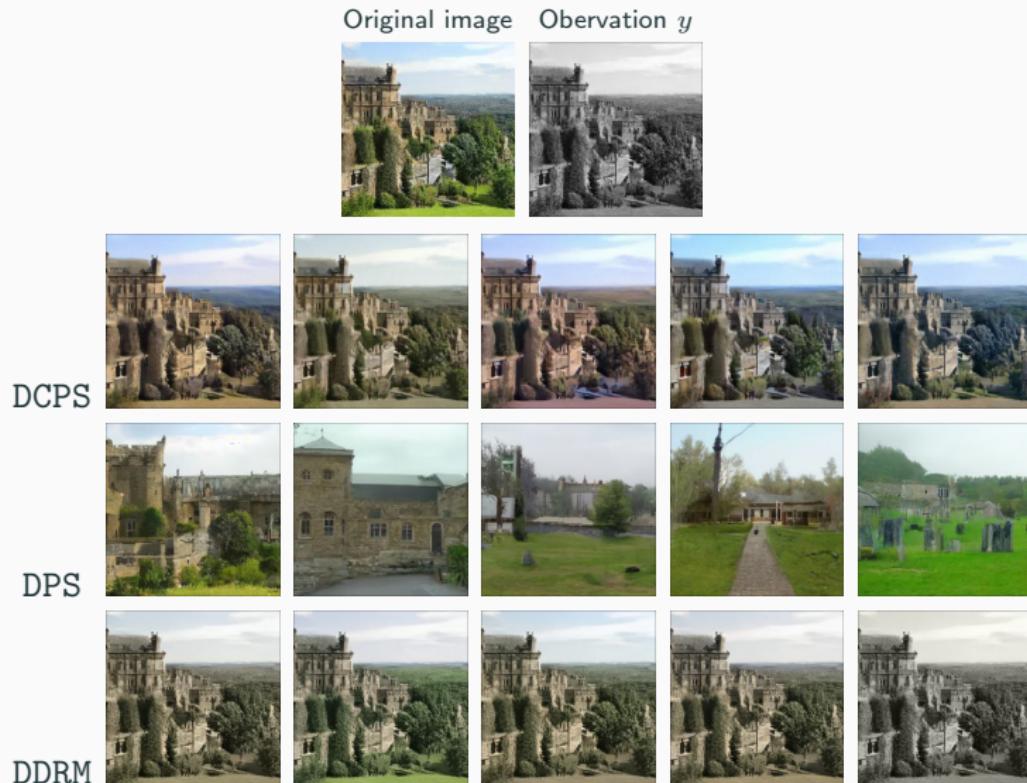
DPS



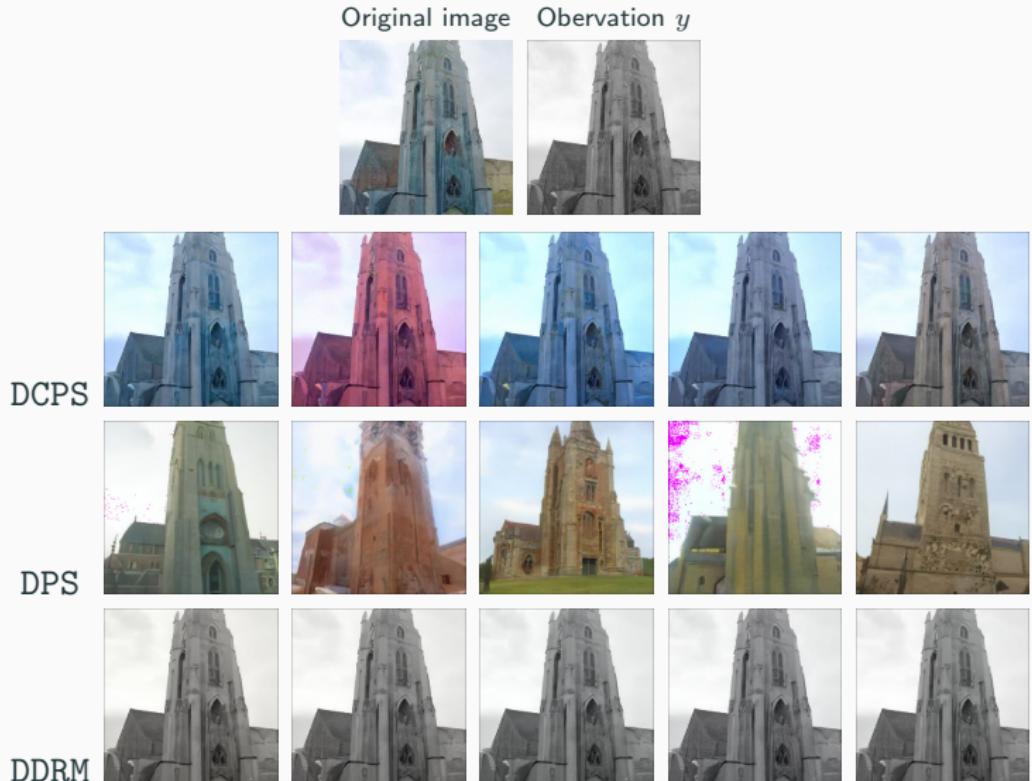
DDRM



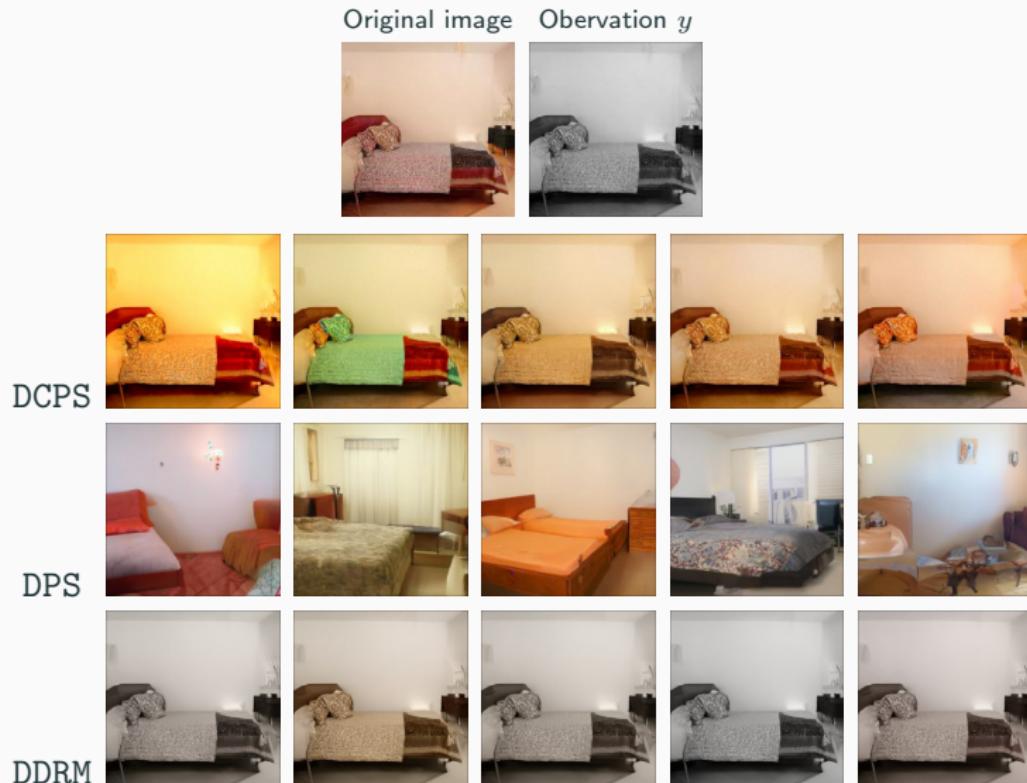
# Colorization experiments



# Colorization experiments



# Colorization experiments



*Thank you!*

# Bibliography i

## References

---

- Ackley, D. H., Hinton, G. E., and Sejnowski, T. J. (1985). A learning algorithm for boltzmann machines. *Cognitive science*, 9(1):147–169.
- Batzolis, G., Stanczuk, J., Schönlieb, C.-B., and Etmann, C. (2021). Conditional image generation with score-based diffusion models. *arXiv preprint arXiv:2111.13606*.
- Bishop, C. M. (2006). *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg.
- Boys, B., Girolami, M., Pidstrigach, J., Reich, S., Mosca, A., and Akyildiz, O. D. (2023). Tweedie moment projected diffusions for inverse problems. *arXiv preprint arXiv:2310.06721*.

## Bibliography ii

- Brosse, N., Durmus, A., Moulines, É., and Sabanis, S. (2019). The tamed unadjusted langevin algorithm. *Stochastic Processes and their Applications*, 129(10):3638–3663.
- Cappe, O., Moulines, E., and Ryden, T. (2005). *Inference in Hidden Markov Models (Springer Series in Statistics)*. Springer-Verlag, Berlin, Heidelberg.
- Cérou, F., Del Moral, P., Furon, T., and Guyader, A. (2012). Sequential Monte Carlo for rare event estimation. *Statistics and computing*, 22(3):795–808.
- Chopin, N., Papaspiliopoulos, O., et al. (2020). *An introduction to sequential Monte Carlo*, volume 4. Springer.
- Chung, H., Kim, J., Mccann, M. T., Klasky, M. L., and Ye, J. C. (2023). Diffusion posterior sampling for general noisy inverse problems. In *The Eleventh International Conference on Learning Representations*.

## Bibliography iii

- Dai, C., Heng, J., Jacob, P. E., and Whiteley, N. (2022). An invitation to sequential monte carlo samplers. *Journal of the American Statistical Association*, 117(539):1587–1600.
- Del Moral, P. (2004). Feynman-kac formulae. In *Feynman-Kac Formulae*, pages 47–93. Springer.
- Del Moral, P. (2013). *Mean Field Simulation for Monte Carlo Integration*. CRC Press.
- Dhariwal, P. and Nichol, A. (2021). Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794.
- Dockhorn, T., Vahdat, A., and Kreis, K. (2022). Score-based generative modeling with critically-damped langevin diffusion. In *International Conference on Learning Representations (ICLR)*.

## Bibliography iv

- Finzi, M. A., Boral, A., Wilson, A. G., Sha, F., and Zepeda-Núñez, L. (2023). User-defined event sampling and uncertainty quantification in diffusion models for physical dynamical systems. In *International Conference on Machine Learning*, pages 10136–10152. PMLR.
- Gelman, A. and Meng, X.-L. (1998). Simulating normalizing constants: From importance sampling to bridge sampling to path sampling. *Statistical Science*, pages 163–185.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680.
- Gordon, N. J., Salmond, D. J., and Smith, A. F. M. (1993). Novel approach to nonlinear/ non-Gaussian Bayesian state estimation. *IEE Proceedings-F*, 140(2):107–113.

## Bibliography v

- Ho, J., Jain, A., and Abbeel, P. (2020). Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851.
- Kawar, B., Elad, M., Ermon, S., and Song, J. (2022). Denoising diffusion restoration models. *Advances in Neural Information Processing Systems*, 35:23593–23606.
- Kingma, D. P. and Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
- Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R., and Van Gool, L. (2022). Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11461–11471.
- Pitt, M. K. and Shephard, N. (1999). Filtering via simulation: Auxiliary particle filters. *J. Amer. Statist. Assoc.*, 94(446):590–599.

## Bibliography vi

- Rezende, D. and Mohamed, S. (2015). Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR.
- Saharia, C., Ho, J., Chan, W., Salimans, T., Fleet, D. J., and Norouzi, M. (2022). Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4713–4726.
- Sohl-Dickstein, J., Weiss, E., Maheswaranathan, N., and Ganguli, S. (2015). Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR.
- Song, J., Meng, C., and Ermon, S. (2021a). Denoising diffusion implicit models. In *International Conference on Learning Representations*.

## Bibliography vii

- Song, J., Vahdat, A., Mardani, M., and Kautz, J. (2023). Pseudoinverse-guided diffusion models for inverse problems. In *International Conference on Learning Representations*.
- Song, Y. and Ermon, S. (2019). Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32.
- Song, Y. and Ermon, S. (2020). Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. (2021b). Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*.

## Bibliography viii

- Syed, S., Romaniello, V., Campbell, T., and Bouchard-Côté, A. (2021). Parallel tempering on optimized paths. In *International Conference on Machine Learning*, pages 10033–10042. PMLR.
- Tashiro, Y., Song, J., Song, Y., and Ermon, S. (2021). CSDI: Conditional score-based diffusion models for probabilistic time series imputation. *Advances in Neural Information Processing Systems*, 34:24804–24816.
- Trippe, B. L., Yim, J., Tischer, D., Baker, D., Broderick, T., Barzilay, R., and Jaakkola, T. S. (2023). Diffusion probabilistic modeling of protein backbones in 3d for the motif-scaffolding problem. In *The Eleventh International Conference on Learning Representations*.
- Wu, L., Trippe, B. L., Naesseth, C. A., Cunningham, J. P., and Blei, D. (2023). Practical and asymptotically exact conditional sampling in diffusion models. In *Thirty-seventh Conference on Neural Information Processing Systems*.

## Bibliography ix

Zhang, G., Ji, J., Zhang, Y., Yu, M., Jaakkola, T., and Chang, S. (2023). Towards coherent image inpainting using denoising diffusion implicit models. *arXiv preprint arXiv:2304.03322*.