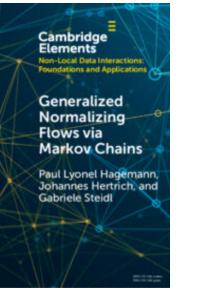


0. Motivation
1. Splitting Algorithms
2. Optimal Transport, Sinkhorn Algorithm and SMART
3. Normalizing Flows
4. Generalized Normalizing Flows

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58



## Lecture 3: Normalizing Flows

Gabriele Steidl

Applied Mathematics - Imaging Sciences

TU Berlin

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Course: PGMO and M2 Optimization

## February 7. & 9. February 2023, École Polytechnique, Paris



Mathematisches  
Forschungsinstitut  
Oberwolfach



## Oberwolfach Seminar

### Variational and Information Flows in Machine Learning and Optimal Transport

Organizers: Wuchen Li, Columbia  
Bernhard Schmitzer, Göttingen  
Gabriele Steidl, Berlin  
Francois-Xavier Vialard, Paris

Date (ID): 19 – 25 November 2023 (2347b)  
Deadline: 1 September 2023

Variational and stochastic flows are now ubiquitous in machine learning and generative modeling. Indeed, many such models can be interpreted as flows from a latent distribution to the sample distribution and training corresponds to finding the right flow vector field. Optimal transport and diffeomorphic flows provide powerful frameworks to analyze such trajectories of distributions with elegant notions from differential geometry, such as geodesics, gradient and Hamiltonian flows. Recently, mean field control and mean field games offer a general optimal control variational problems on the learning problem. How do these tools lead us to a better understanding and further development of machine learning and generative models?

The Oberwolfach Seminar will address the topic from different points of view taking in particular recent developments in machine learning into account. The target audience is PhD students and post-doctoral researchers wishing to be quickly immersed in this modern, active research area. Priority will be given to young, motivated researchers.

Please see the website of the seminar for detailed information:

[www.mfo.de/occasion/2347b](http://www.mfo.de/occasion/2347b)

The seminar takes place at the Mathematisches Forschungsinstitut Oberwolfach. The Institute covers board and lodging. By the support of the Carl Friedrich von Siemens Foundation travel expenses can be reimbursed up to 150 EUR in average per person (against copies of travel receipts). The number of participants is restricted to 25.

**Applications including title, ID and date** of the intended seminar, together with **one pdf-file attached** containing

- full name and address, incl. e-mail address
- short CV and publication list
- present position, university
- name of supervisor of Ph.D. thesis
- a short summary of previous work and interest

should be **sent by e-mail** via [seminars@mfo.de](mailto:seminars@mfo.de) until 1 September 2023 to:

Prof. Dr. Matthias Hieber  
Mathematisches Forschungsinstitut Oberwolfach  
Schwarzwaldstr. 9 – 11  
77709 Oberwolfach  
Germany



[www.mfo.de/scientific-program/meetings/oberwolfach-seminars](http://www.mfo.de/scientific-program/meetings/oberwolfach-seminars)

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

- ◆ Generative Adversarial Networks (GANs) (Goodfellow et al. 2014)
- ◆ Variational Auto-Encoders (Kingma/Welling 2014)
- ◆ Diffusion Flows (Zhang/Chen 2021, ...  )
- ◆ Invertible Networks
  - Residual Networks (Behrmann, Chen et al. 2019)
  - Normalizing Flows - Directly Invertible Neural Networks (Dingh et al. 2017, Aridizzone et al. 2019)
    - for continuous normalizing flows  
see, e.g. Ruthotto/Haber 2020, Hagemann/Hertrich/St. Overview paper:  
Generalized Normalizing Flows via Markov Chains 2021

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Push Forward Measures

- ◆  $T : X \rightarrow Y$  (Borel) measurable,  $\mu \in \mathcal{P}(X)$

Push forward measure

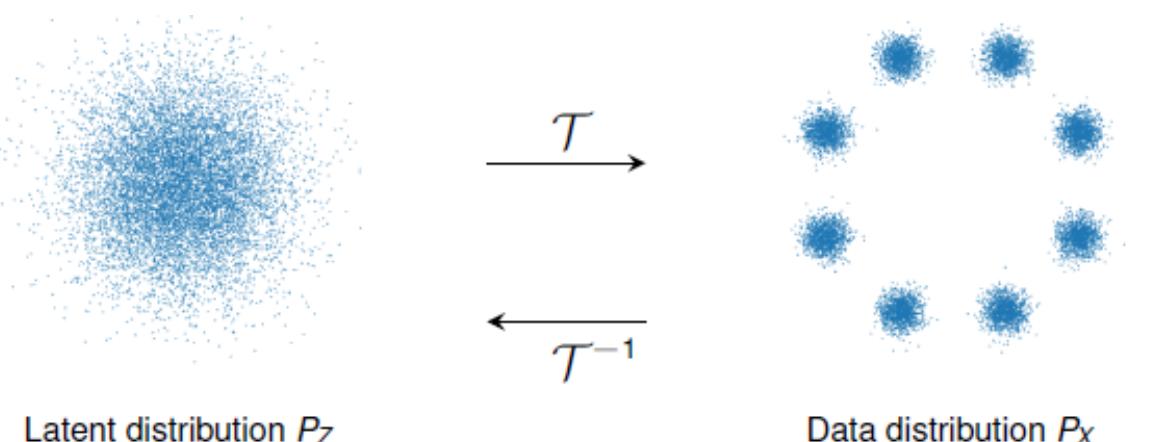
$$T_{\#}\mu := \mu \circ T^{-1}$$

Relation

$$\int_{T^{-1}(A)} h(T(x)) d\mu(x) = \int_A h(y) d\underbrace{(T_{\#}\mu)}_{\nu}(y)$$

Change of variable formula: in case of existing density  $p_\mu$  and a diffeomorphism  $T$ :

$$p_{T_{\#}\mu}(y) = p_\mu(T^{-1}(y)) |\det \nabla T^{-1}(y)|$$



1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Outline

1. Residual Networks (ResNet)
  - Proximal NNs within ResNets
2. Normalizing Flows
  - Applications in Inverse Problems - The Power of Patches

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Residual Networks (ResNets)

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

**Residual Networks:** composition  $\mathcal{T}_N \circ \dots \circ \mathcal{T}_1$  of layers of the form

$$\mathcal{T}_k(x) = x + \Phi(x; \theta_k), \quad \Phi(x; \theta_k) \text{ subnetworks}$$

All  $\mathcal{T}_k$  are networks themselves. But they are chosen to be invertible as follows.

**Invertible Residual Networks:**  $x^{(0)} = y (= x + \Phi(x; \theta_k))$

$$x^{(r+1)} = y - \Phi(x^{(r)}; \theta_k).$$

Sequence  $(x^{(r)})_r$  converges to  $\mathcal{T}_k^{-1}(y)$  if  $\text{Lip}(\Phi(\cdot; \theta_k)) < 1$

BIG effort to ensure Lipschitz continuity in learning process!

**Proposal:**

$$\mathcal{T}_k(x) = x + \gamma_k \underbrace{\Phi(x; \theta_k)}_{PNN}, \quad \gamma_k > 0$$

Refs: Behrmann et al. 2019, Chen et al. 2019, J. Hertrich: Proximal Residual Flows for Bayesian Inverse Problems, 2022

# Residual Networks (ResNets)

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

**Theorem** Let  $\Phi = (1 - t)I + tR$ ,  $t \in (0, 1)$  be an averaged operator. If  $\frac{1}{2} < t \leq 1$ , let  $0 < \gamma < \frac{1}{2t-1}$  and  $\gamma > 0$  otherwise. Then  $\mathcal{T}(x) := x + \gamma\Phi(x)$  is invertible and  $\mathcal{T}^{-1}(y)$  is the limit of

$$x^{(r+1)} = \frac{1}{1 + \gamma - \gamma t}y - \frac{\gamma t}{1 + \gamma - \gamma t}R(x^{(r)}),$$

Proof:  $R$  is 1-Lipschitz. For  $\frac{1}{2} < t \leq 1$  we have  $\gamma < \frac{1}{2t-1} \Leftrightarrow \frac{\gamma t}{1+\gamma-\gamma t} < 1$ . By Banach's fixed point theorem the series converges to the unique fixed point of

$$x = \frac{1}{1 + \gamma - \gamma t}y - \frac{\gamma t}{1 + \gamma - \gamma t}R(x)$$

Then

$$\begin{aligned} y &= (1 + \gamma - \gamma t)x + \gamma tR(x) \\ &= x + \gamma((1 - t)x + tR(x)) \\ &= x + \gamma\Phi(x) \end{aligned}$$

Thus,  $x = \mathcal{T}^{-1}(y)$ .

# Learning Residual Networks (ResNets)

Minimization of log likelihood loss:

$$\mathcal{L}(\theta) = \sum_{i=1}^N \log p_{\mathcal{T}_\theta^{-1} P_Z}(x_i).$$

We have to evaluate and differentiate  $\log |\det \nabla \mathcal{T}|$  which is computationally intractable in high dimensions.

**Theorem** Let  $Q$  be a random variable on  $\mathbb{Z}_{>0}$  such that  $P(Q = k) > 0$  for all  $k \in \mathbb{Z}_{>0}$  and  $p_k := P(Q \geq k)$ . Let  $\mathcal{T}(x) = x + \Phi(x)$ , where  $\Phi: \mathbb{R}^n \rightarrow \mathbb{R}^n$  is differentiable and fulfills  $\text{Lip}(\Phi) < 1$ . Then, it holds

$$\log(|\nabla \mathcal{T}(x)|) = \mathbb{E}_{v \sim \mathcal{N}(0, I), q \sim P_Q} \left[ \sum_{k=1}^q \frac{(-1)^{k+1}}{k} \frac{v^\top (\nabla \Phi(x))^k v}{p_k} \right]$$

and

$$\frac{\partial}{\partial \theta} \log(|\nabla \mathcal{T}(x)|) = \mathbb{E}_{v \sim \mathcal{N}(0, I), q \sim P_Q} \left[ \left( \sum_{k=0}^q \frac{(-1)^k}{p_k} v^\top (\nabla \Phi(x))^k \right) \frac{\partial (\nabla \Phi(x))}{\partial \theta} v \right].$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Conditional Proximal Residual Flows

**Setting:**  $Y = F(X) + \eta$ ,

where  $F: \mathbb{R}^n \rightarrow \mathbb{R}^d$  is an ill-posed/ill-conditioned forward operator and  $\eta$  is some noise.

**Aim:** NN for reconstructing all posterior distributions  $P_{X|Y=y}$ ,  $y \in \mathbb{R}^d$   
 $\mathcal{T}_\theta: \mathbb{R}^d \times \mathbb{R}^n \rightarrow \mathbb{R}^n$  such that  $\mathcal{T}(y, \cdot)$  is invertible for all  $y \in \mathbb{R}^d$  and

$$P_{X|Y=y} = \mathcal{T}_\theta(y, \cdot)_\#^{-1} P_Z.$$

Learn  $\mathcal{T}_\theta$  from i.i.d. samples  $(x_1, y_1), \dots, (x_N, y_N)$  of  $(X, Y)$  using the maximum likelihood loss

$$\mathcal{L}(\theta) = \sum_{i=1}^N p_{\mathcal{T}_\theta(y_i, \cdot)_\#^{-1} P_Z}(x_i) \approx \mathbb{E}_{y \sim P_Y} [\text{KL}(P_{X|Y=y}, P_{\mathcal{T}_\theta(y, \cdot)_\#^{-1} P_Z})] + \text{const.}$$

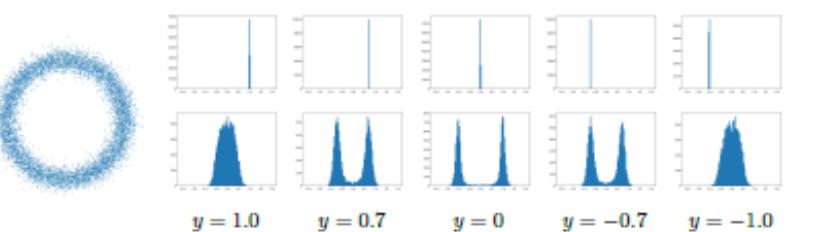


Figure 2: Left: Samples from the prior distribution of  $X$  for the circle example. Right: Histograms of samples from the reconstructed posterior distribution  $P_{X|Y=y} \approx \mathcal{T}(\cdot, y)_\#^{-1} P_Z$  for  $y \in \{1, 0.7, 0, -0.7, -1\}$  within the circle example. Top: first coordinate, Bottom: second coordinate.

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

## Remark

$$J(\theta) = \mathbb{E}_{x \sim P} [\text{trace}(\nabla_x^2 \log(q(x, \theta))) + \frac{1}{2} \|\nabla_x \log(q(x, \theta))\|^2].$$

Lemma:

Let  $A \in \mathbb{R}^{d \times d}$  be an arbitrary matrix and let  $V$  be a random variable with  $\mathbb{E}(V) = 0$  and  $\text{Cov}(V) = I$ . Then, it holds

$$\text{trace}(A) = \mathbb{E}(V^\top A V).$$

Using the lemma,  $J$  can be rewritten as

$$J(\theta) = \mathbb{E}_{x \sim P, v \sim \mathcal{N}(0, 1)} [v^\top \nabla_x^2 \log(q(x, \theta)) v + \frac{1}{2} \|\nabla_x \log(q(x, \theta))\|^2].$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Normalizing Flows (NFs)

Parameter depending concatenation of diffeomorphisms of the form (up to some permutation matrices)

$$\mathcal{T}(\cdot; \theta) = \mathcal{T}_T \circ \cdots \circ \mathcal{T}_1$$

Invertibility by special structure of  $\mathcal{T}_k$ ,  $k = 1, \dots, T$ :

1. Real NVP (real-valued volume-preserving transformations) (Dingh et al. 2017)

$$\mathcal{T}_k(z_1, z_2) = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} := \begin{pmatrix} z_1 \\ z_2 e^{\mathbf{s}_k(z_1)} + \mathbf{t}_k(z_1) \end{pmatrix}$$

$$\mathcal{T}_k^{-1}(x_1, x_2) = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} x_1 \\ (x_2 - \mathbf{t}_k(x_1)) e^{-\mathbf{s}_k(x_1)} \end{pmatrix}$$

with **neural networks**  $s_k, t_k$ ,  $k = 1, \dots, T$ ,  $x_j, z_j \in \mathbb{R}^{n_j}$ ,  $j = 1, 2$ .

2. Other architecture (Aridizzone et al. 2019):

$$\mathcal{T}_k(z_1, z_2) = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} := \begin{pmatrix} z_1 e^{\mathbf{s}_{k,2}(z_2)} + \mathbf{t}_{k,2}(z_2) \\ z_2 e^{\mathbf{s}_{k,1}(x_1)} + \mathbf{t}_{k,1}(x_1) \end{pmatrix}$$

$$\mathcal{T}_k^{-1}(x_1, x_2) = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} (x_1 - \mathbf{t}_{k,2}(z_2)) e^{-\mathbf{s}_{k,2}(z_2)}, \\ (x_2 - \mathbf{t}_{k,1}(x_1)) e^{-\mathbf{s}_{k,1}(x_1)} \end{pmatrix}$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Normalizing Flows (NFs)

Since  $T_k = T_{2,k} \circ T_{1,k}$  with

$$T_{1,k}(z_1, z_2) = \begin{pmatrix} x_1 \\ z_2 \end{pmatrix} := \begin{pmatrix} z_1 e^{s_{k,2}(z_2)} + t_{k,2}(z_2) \\ z_2 \end{pmatrix},$$

$$T_{2,k}(x_1, z_2) = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} := \begin{pmatrix} x_1 \\ z_2 e^{s_{k,1}(x_1)} + t_{k,1}(x_1) \end{pmatrix},$$

we have

$$\nabla T_{1,k}(z_1, z_2) = \begin{pmatrix} \text{diag}\left(e^{s_{k,2}(z_2)}\right) & \text{diag}\left(\nabla_{z_2}\left(z_1 e^{s_{k,2}(z_2)} + t_{k,2}(z_2)\right)\right) \\ 0 & I_{d_2}, \end{pmatrix}$$

so that  $\det \nabla T_{1,k}(z_1, z_2) = \prod_{k=1}^{d_1} e^{(s_{k,2}(z_2))_k}$  and similarly for  $\nabla T_{2,k}$ . Applying the chain rule and noting that  $\det(AB) = \det(A)\det(B)$ ,

$$\log(|\det(\nabla \mathcal{T}(z))|) = \sum_{k=1}^T (\text{sum}\left(s_{k,2}\left((z^k)_2\right)\right) + \text{sum}\left(s_{k,1}\left((T_{1,k}z^k)_1\right)\right)),$$

where  $\text{sum}$  denotes the sum of the components of the respective vector,  $z^1 := z$  and  $z^k = T_{k-1}z^{k-1}$ ,  $k = 2, \dots, T$ .

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Learning Normalizing Flows

**Aim:**  $P_X \approx \mathcal{T}_\# P_Z := P_Z \circ \mathcal{T}^{-1}$

**Loss Function:** Kullback-Leibler divergence (KL)

$$\begin{aligned}\mathcal{L}_{\text{NF}}(\theta) &= \text{KL}(P_X, \mathcal{T}_\# P_Z) = \int \log \left( \frac{p_X(x)}{p_{\mathcal{T}_\# P_Z}(x)} \right) p_X(x) dx \\ &= \underbrace{\int \log(p_X(x)) p_X(x) dx}_{\text{const}} - \int \log(p_{\mathcal{T}_\# P_Z}(x)) p_X(x) dx\end{aligned}$$

with transformation formula  $p_{\mathcal{T}_\# P_Z} = p_Z(\mathcal{T}^{-1}) |\det \nabla \mathcal{T}^{-1}|$  and Gaussian distribution  $P_Z$

$$\begin{aligned}\mathcal{L}_{\text{NF}}(\theta) &\sim - \int \log(p_Z(\mathcal{T}^{-1}(x)) |\det \nabla \mathcal{T}^{-1}(x)|) p_X(x) dx \\ &= -\mathbb{E}_{x \sim P_X} [\log p_Z(\mathcal{T}^{-1})] - \mathbb{E}_{x \sim P_X} [\log(|\det \nabla \mathcal{T}^{-1}|)] \\ &= \|\mathcal{T}^{-1}(\cdot)\|_{L_2(dP_X)}^2 - \mathbb{E}_{x \sim P_X} [(\log |\det \nabla \mathcal{T}^{-1}|)]\end{aligned}$$

Refs: SGD (Optimizer Adam: Kingma et al. 2015), Inertial Stoch. PALM (Hertrich/St. 2022)

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

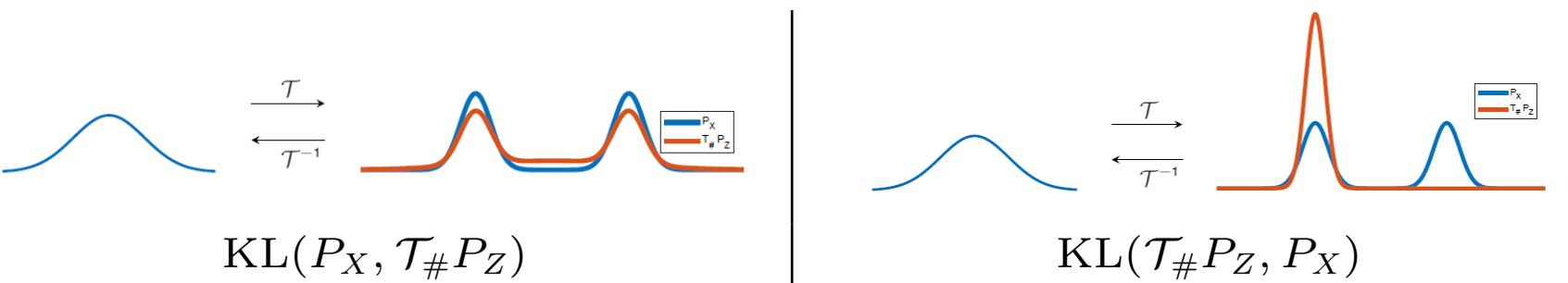
# Remark on Kullback-Leibler Divergence

KL is (only) the Bregman distance of the Shannon entropy:

- $\text{KL}(\mu, \nu) \geq 0$  for all  $\mu, \nu \in \mathcal{P}(\mathbb{R}^d)$  with equality  $\Leftrightarrow$  if  $\nu = \mu$
- **not symmetric** and no triangular inequality
- finite if  $\mu \ll \nu$  and  $\text{KL}(\mu, \nu) = +\infty$  otherwise

Different properties of  $\text{KL}(P_X, \mathcal{T}_\# P_Z)$  (forward KL) and  $\text{KL}(\mathcal{T}_\# P_Z, P_X)$  (backward KL):

- ◆ inverse problems: operator known **versus** operator not known
- ◆ mode covering (unrealistic samples possible) **versus** mode seeking (mode collapse)



Refs: Backward KL: Kruse et al. 2020, Sun/Bouman 2021, Altekrüger/Hertrich: WPPFlows 2022,

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Applications in Inverse Problems - The Power of Patches

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

Inverse Problem for a certain class of images:

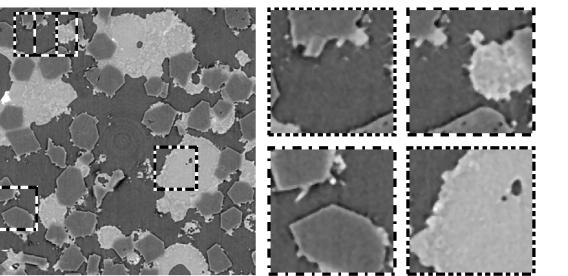
$$y = \text{noisy}(F(x))$$

Find solution as minimizer of variational model:

$$\mathcal{J}(x) = \underbrace{\mathcal{D}_F(x, y)}_{\text{data term}} + \lambda \underbrace{\mathcal{R}(x)}_{\text{regularizer}}, \quad \lambda > 0$$

Idea: Learn regularizer from many patches  $P_i(x_j)$ ,  $i = 1, \dots, N$  of few images

$x_j$ ,  $j = 1, \dots, n$



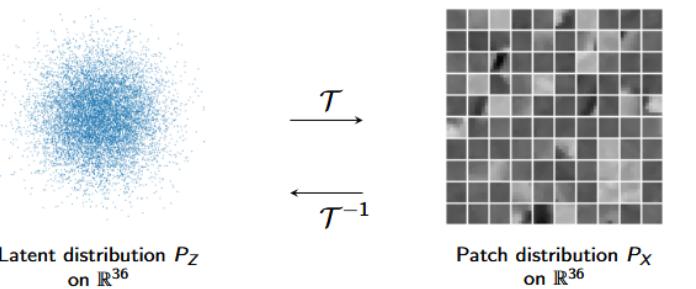
The Power of Patches

Refs: Buardes et al. 2005, Dabov et al. (BM3D) 2008; Lebrun/Morel (Denoising cuisine) 2013; Bortoli/Desolneux/Galerne/Leclaire 2019 ...., Laus et al. 2017, Houdard et al. 2018; Hertrich et al. 2021 ...

# Applications in Inverse Problems

## 1. Learn Normalizing Flow $\mathcal{T} = \mathcal{T}(\theta, \cdot)$

$$\mathcal{L}_{\text{NF}}(\theta) = \sum_{j=1}^n \sum_{i=1}^N \frac{\|\mathcal{T}^{-1}(P_i(\mathbf{x}_j))\|^2}{2} - \log |\det \nabla \mathcal{T}^{-1}(P_i(\mathbf{x}_j))|$$



## 2. Variational Model

$$\mathcal{J}(x) = \mathcal{D}_F(x, y) + \lambda \text{patchNR}(x)$$

negative log likelihood of all patches under the probability distribution learned by the patchNR  $\mathcal{T} = \mathcal{T}(\theta, \cdot)$

$$\text{patchNR}(x) := \left( \frac{1}{N} \sum_{i=1}^N \frac{\|\mathcal{T}^{-1}(P_i(x))\|^2}{2} - \log |\det \nabla \mathcal{T}^{-1}(P_i(x))| \right)$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Applications in Inverse Problems

Bayesian Approach - MAP:

$$\begin{aligned}
 \hat{x} &\in \operatorname{argmax}_x \log(p_{X|Y=y}(x)) \\
 &= \operatorname{argmin}_x \left\{ -\log(p_{Y|X=x}(y)) - \log(p_X(x)) \right\} \\
 &= \operatorname{argmin}_x \left\{ \underbrace{\mathcal{D}_F(x, y)}_{\text{data-fidelity term}} + \lambda \underbrace{\mathcal{R}(x)}_{\text{regularizer}} \right\}
 \end{aligned}$$

**Proposition:**

Let  $P_Z = \mathcal{N}(0, I)$  and let  $\mathcal{T}: \mathbb{R}^s \rightarrow \mathbb{R}^s$  be a bi-Lipschitz diffeomorphism. Then

$$\exp(-\lambda \operatorname{patchNR}(x)) \in L^1(\mathbb{R}^d), \quad \lambda > 0.$$

**Advantages:**

- ◆ **Few** training data: often not many training images are available - need just patches! Save electrical energy.
- ◆ **Flexibility** of operator and image classes: regularizer fits for every operator in the corresponding image classes

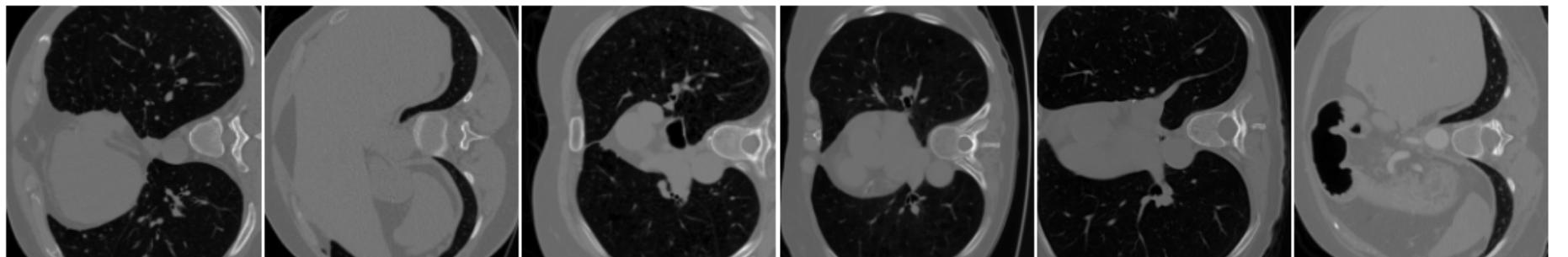
1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Computerized Tomography

Data term: Poisson (like) noise

$$\mathcal{D}(Ax, y) = \sum_{i=1}^d e^{-(Ax)_i} N_0 - e^{-y_i} N_0 (- (Ax)_i + \log(N_0))$$

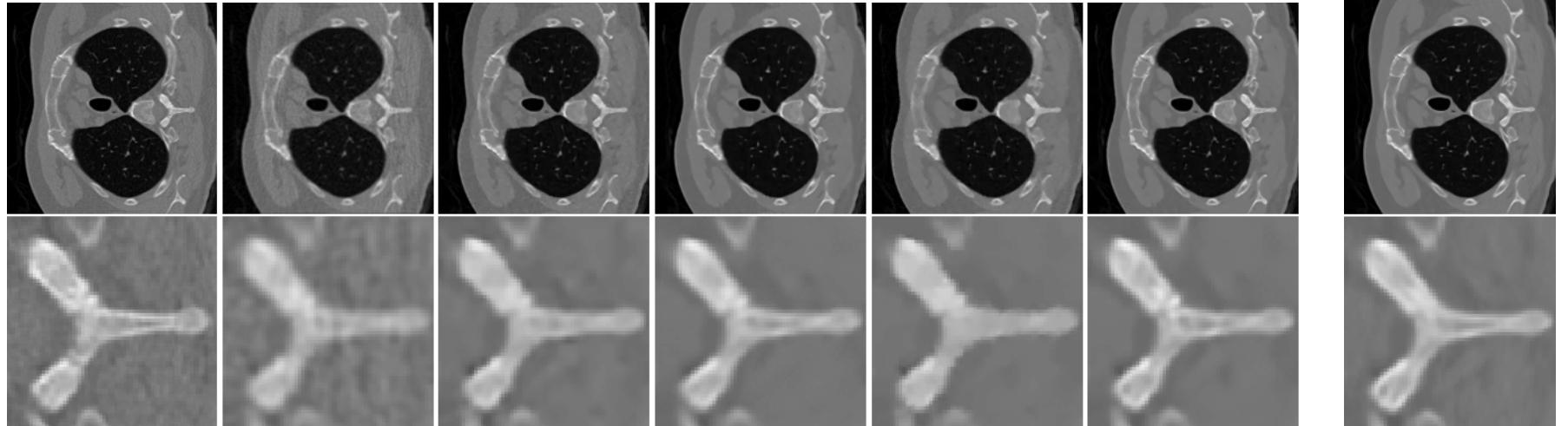
Regularizer:  $\mathcal{R}(x) = \text{patchNR}(x)$  learned from  $n = 6$  images



1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Low Dose Computerized Tomography

Results: Low Dose Tomography (full angle)



Ground truth

FBP

DIP+TV

EPLL

localAR

patchNR

FBP+UNet

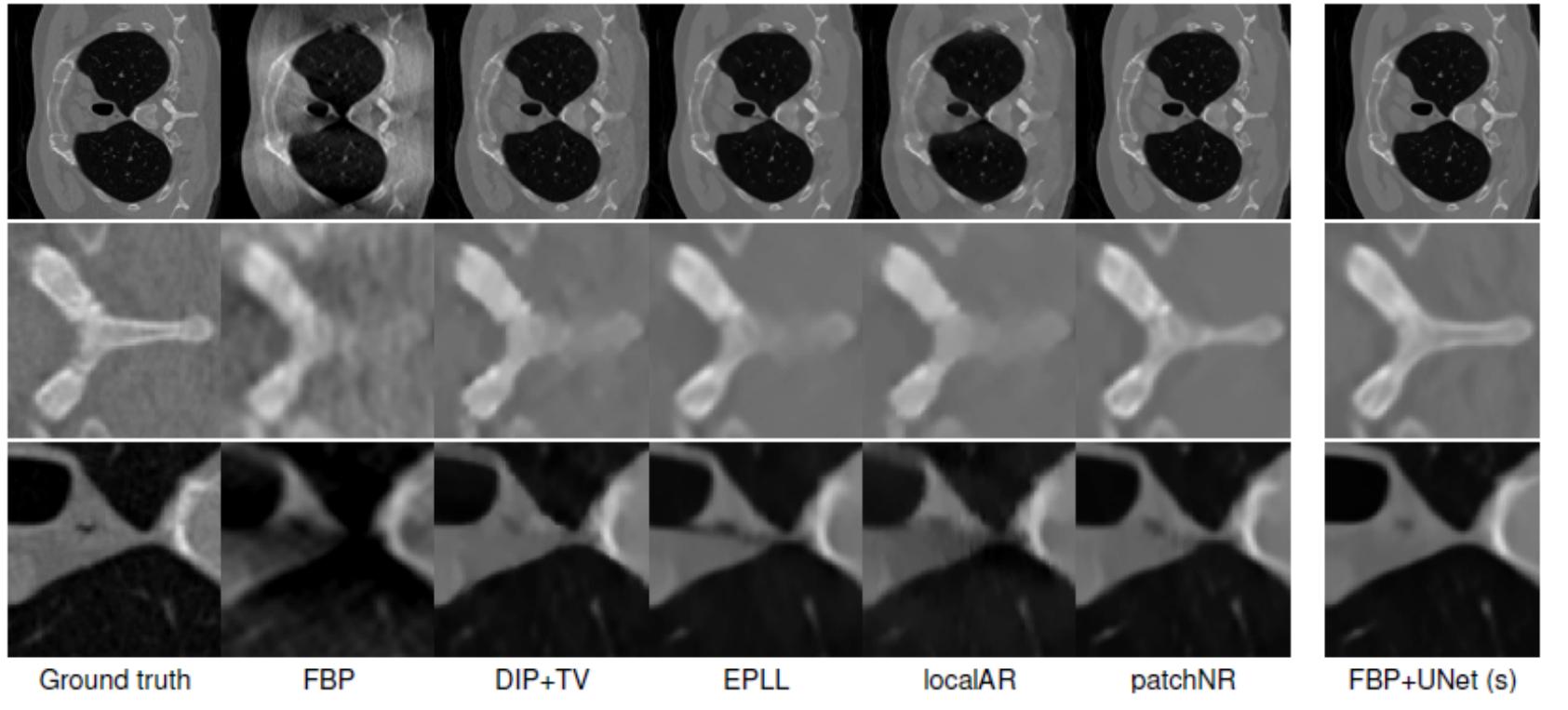
1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

	FBP	DIP + TV	EPLL	localAR	patchNR	FBP+UNet (data-based)
PSNR	$30.37 \pm 2.95$	$34.45 \pm 4.20$	$34.89 \pm 4.41$	$33.64 \pm 3.74$	$35.19 \pm 4.52$	$35.48 \pm 4.52$
SSIM	$0.739 \pm 0.141$	$0.821 \pm 0.147$	$0.821 \pm 0.154$	$0.807 \pm 0.145$	$0.829 \pm 0.152$	$0.837 \pm 0.143$
Runtime	0.03s	1514.33s	36.65s	30.03s	48.39s	0.46s

Refs: DIP-TV: Ulyanov et al. 2018, EPLL: Zoran et al. 2011, AR: Lunz et al. 2018, LocalAR: Prost et al. 2021, FBP+UNet: Jin et al. 2017, **35820 supervised samples**

# Limited Angle Computerized Tomography

Results: Limited Angles -  $36^\circ/180^\circ$  (same  $\mathcal{R}$  for Low Dose Tomography!)



	FBP	DIP + TV	EPLL	localAR	patchNR	FBP+UNet (data-based)
PSNR	$21.96 \pm 2.25$	$32.57 \pm 3.25$	$32.78 \pm 3.46$	$31.06 \pm 2.95$	$33.20 \pm 3.55$	$33.75 \pm 3.58$
SSIM	$0.531 \pm 0.097$	$0.803 \pm 0.146$	$0.801 \pm 0.151$	$0.779 \pm 0.142$	$0.811 \pm 0.151$	$0.820 \pm 0.140$
Runtime	0.02s	1770.89s	127.21s	53.47s	485.93s	0.53s

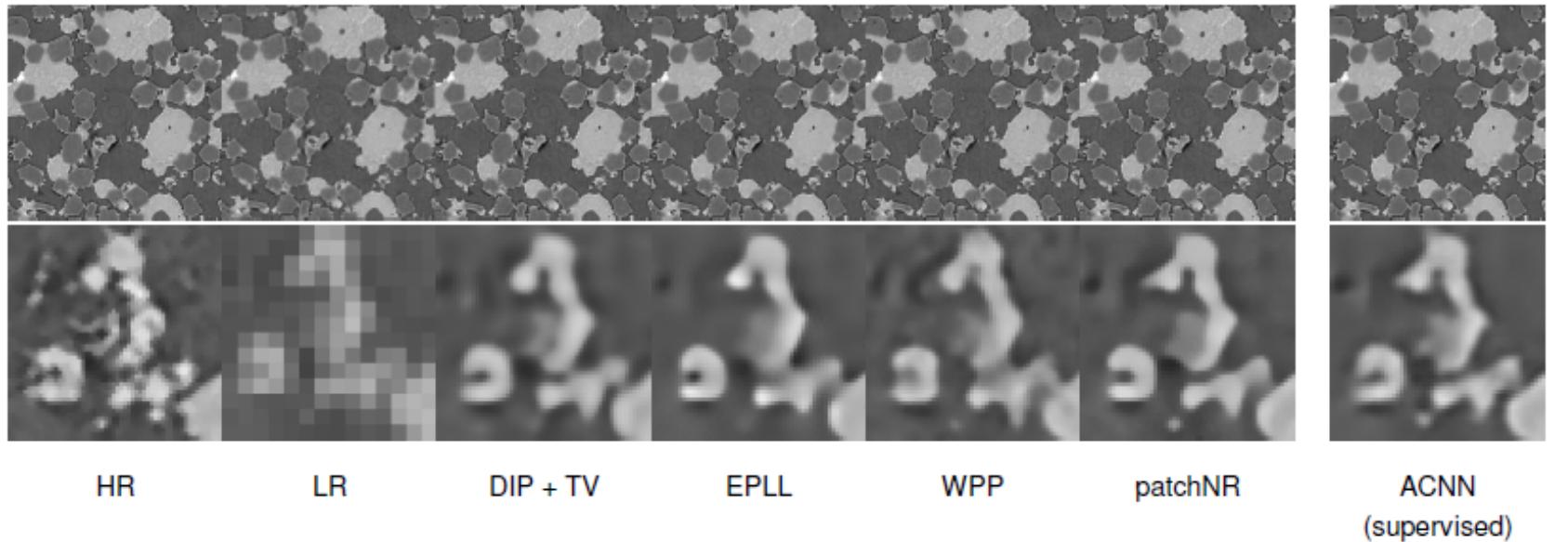
1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Superresolution

Data term: Gaussian noise  $\mathcal{D}(Ax, y) = \frac{1}{2}\|Ax - y\|^2$  stride: 4

Regularizer: learned from 1 high resolution synchotron image

Results:



	bicubic (not shown)	DPIR (not shown)	DIP + TV	EPLL	WPP	patchNR	ACNN (data-based)
PSNR	$25.63 \pm 0.56$	$27.78 \pm 0.53$	$27.99 \pm 0.54$	$28.11 \pm 0.55$	$27.80 \pm 0.37$	<b><math>28.53 \pm 0.49</math></b>	$28.89 \pm 0.53$
SSIM	$0.699 \pm 0.012$	$0.770 \pm 0.011$	$0.764 \pm 0.007$	$0.779 \pm 0.010$	$0.749 \pm 0.011$	<b><math>0.780 \pm 0.008</math></b>	$0.804 \pm 0.010$
Runtime	0.0002s	56.62s	234.00s	60.28s	387.28s	150.79s	0.03s

Data: D. Bernard (U Bordeaux), Y. Berthoumieu, JF Aujol (ANR-DFG project)

Refs: WPP: Hertrich et al. 2021, ACNN: Tian et al. 2021

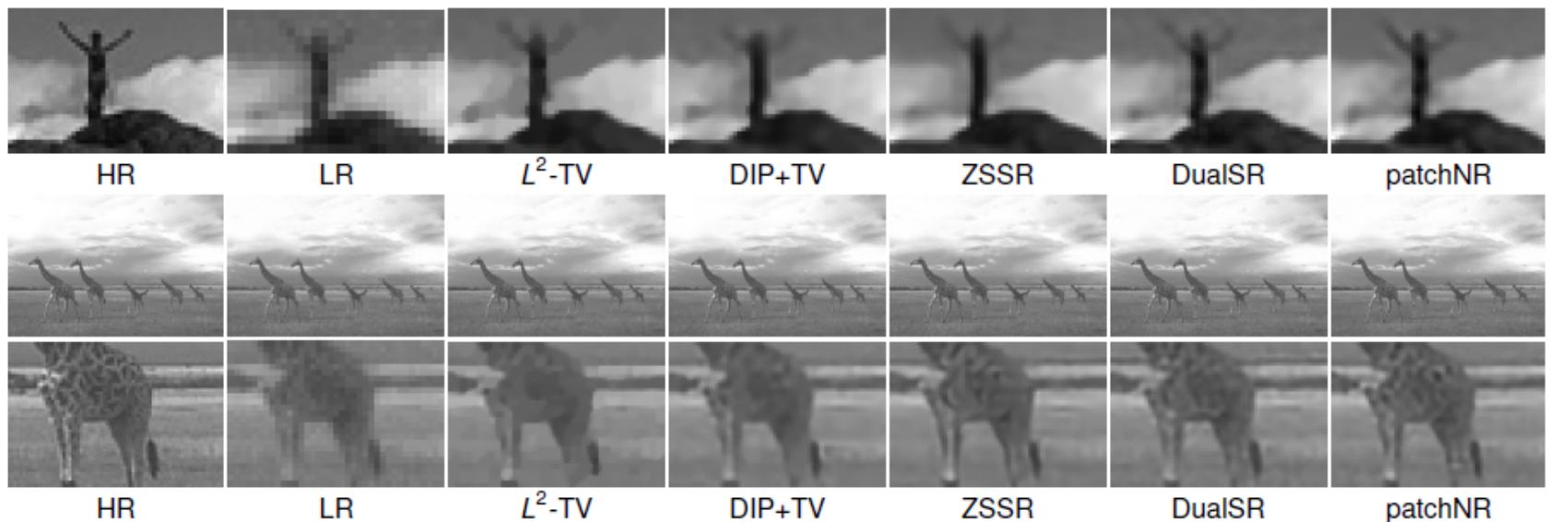
1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Single-Shot Superresolution

**Question:** What to do when we do not have any high-resolution image for training the patchNR?

**Assumption:** Patch distribution of **natural images** is self-similar across the scales

**Regularizer:** learned from 1 low resolution image



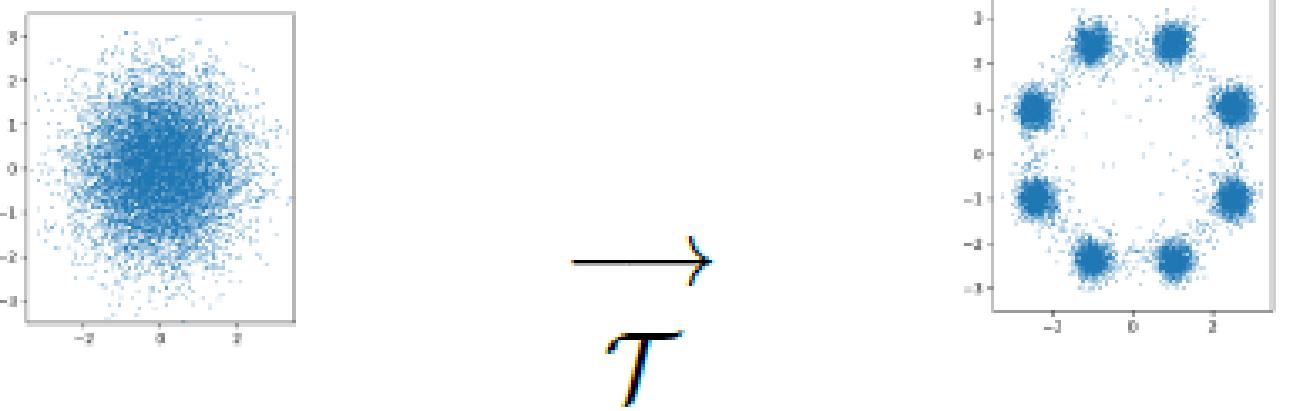
	$L^2$ -TV	DIP+TV	ZSSR	DualSR	patchNR
PSNR	$28.35 \pm 3.55$	$28.44 \pm 3.69$	$28.83 \pm 3.57$	$28.64 \pm 3.47$	<b><math>29.08 \pm 3.58</math></b>
SSIM	$0.820 \pm 0.072$	$0.821 \pm 0.087$	$0.834 \pm 0.066$	$0.829 \pm 0.061$	<b><math>0.846 \pm 0.061</math></b>
Runtime	13.12s	171.51s	56.64s	53.47s	132.36s

Refs: Glasner et al. 2009, ZSSR: Shocher et al. 2018 , Dual SR: Emad et al. 2021

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

## Drawback of NF

- ◆ **Drawback: Lack in expressiveness!**: unimodal distributions are hard to map to **multimodal** and heavy tailed ones
- ◆ Normalizing flows mapping unimodal distributions to multimodal ones must have an **exploding Lipschitz constant!**  
 (Refs Lipschitz: Nagarajan et al. 2018, Gulrajani et al. 2018, Hagemann/Neumayer 2021; Stéphanovitch et al. 2022; Salmona et al. 2022)



- ◆ **(One) Solution:** Application of stochastic steps within a **unifying framework of Markov chains**  
 (Refs Stochastic NFs without Markov chains: Wu/Köhler/Noé 2020, Nielsen/Welling et al. 2021)

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

## Lecture 4: Generalized Normalizing Flows

Gabriele Steidl

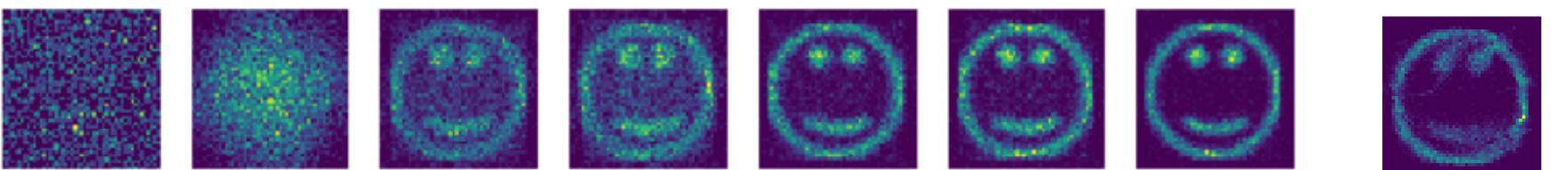
Applied Mathematics - Imaging Sciences  
TU Berlin

# Generalized Normalizing Flows

Markov chain framework includes

- ◆ Normalizing flows
- ◆ Metropolis-Hastings (MH, MCMC) layer
- ◆ Langevin layer
- ◆ VAE layer
- ◆ Diffusion flow layer

**Up to now the only mathematically sound way to combine these layers!**



Alternation of INNs with MCMC layers, starting with the first one.

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Outline

1. Markov Kernels and Markov Chains
2. Normalizing Flows via Markov Chains
3. Generalized Normalizing Flows (GNFs)
4. Conditional GNFs

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Markov Kernels and Markov Chains

- ◆ A **Markov kernel**  $\mathcal{K}: \mathbb{R}^d \times \mathcal{B}(\mathbb{R}^d) \rightarrow [0, 1]$  is a mapping such that
  - i)  $\mathcal{K}(\cdot, B)$  is measurable for any  $B \in \mathcal{B}(\mathbb{R}^d)$ , and
  - ii)  $\mathcal{K}(x, \cdot)$  is a probability measure for any  $x \in \mathbb{R}^d$ .
- ◆ For  $\mu$  on  $\mathcal{P}(\mathbb{R}^n)$ , the measure  $\mu \times \mathcal{K}$  on  $\mathbb{R}^n \times \mathbb{R}^d$  is defined by

$$(\mu \times \mathcal{K})(A \times B) := \int_A \mathcal{K}(x, B) d\mu(x).$$

Definition captures all sets in  $\mathcal{B}(\mathbb{R}^n \times \mathbb{R}^d)$  since the measurable rectangles form a  $\cap$ -stable generator of  $\mathcal{B}(\mathbb{R}^n \times \mathbb{R}^d)$ .

- ◆ For all integrable  $f$  that

$$\int_{\mathbb{R}^n \times \mathbb{R}^d} f(x, y) d(\mu \times \mathcal{K})(x, y) = \int_{\mathbb{R}^n} \int_{\mathbb{R}^d} f(x, y) d\mathcal{K}(x, \cdot)(y) d\mu(x).$$

- ◆ Regular conditional distribution of  $X$  given  $Y$ :  $P_Y$ -a.s. unique Markov kernel  $P_{Y|X=\cdot}(\cdot)$  with

$$P_X \times \underbrace{P_{Y|X=\cdot}(\cdot)}_{\mathcal{K}(\cdot, \cdot)} = P_{(X,Y)}$$

$$\int_{\mathbb{R}^n \times \mathbb{R}^d} f(x, y) dP_{(X,Y)}(x, y) = \int_{\mathbb{R}^n} \int_{\mathbb{R}^d} f(x, y) dP_{Y|X=x}(y) dP_X(x).$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# 1. Markov Kernels and Markov Chains

- ◆  $(X_0, \dots, X_T)$  is called a **Markov chain**, if there exist Markov kernels

$$\mathcal{K}_t = P_{X_t|X_{t-1}=\cdot}(\cdot) : \mathbb{R}^d \times \mathcal{B}(\mathbb{R}^d) \rightarrow [0, 1]$$

such that

$$P_{(X_0, \dots, X_T)} = P_{X_0} \times \mathcal{K}_1 \times \cdots \times \mathcal{K}_T.$$

- ◆ Markov kernels  $\mathcal{K}_t$  are called **transition kernels**
- ◆ If  $\mathcal{K}_t(x, \cdot) = P_{X_t|X_{t-1}=x}$  has a density  $k_t(x, y)$ , and  $P_{X_{t-1}}$  resp.  $P_{X_t}$  have densities  $p_{X_{t-1}}$  resp.  $p_{X_t}$ , then

$$p_{X_t}(y) = \int_{\mathbb{R}^{d_{t-1}}} k_t(x, y) p_{X_{t-1}}(x) dx.$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Example in OT

Discrete Measures

	$v_1 \dots v_n$
$\mu_i$	
$\vdots$	
$\mu_n$	$\pi_{ij}^i$

$$\pi = \sum_{ij} \pi_{ij}^i \delta_{(x_i, y_j)}$$

$$\pi 1 = \mu$$

$$\pi^T 1 = 0$$

$$\mu \times K(x_i, \cdot) = \pi_{ij}^i$$

$$\mu \times K = \pi$$

Markov kernel

$$K(x_i, \cdot) = \sum_j k_{ij} \delta_{y_j}$$

$$P_{Y|X=x_i} \rightarrow P_{Y|X=x_i}(y_j) = k_{ij}$$

$$K^T \mu = 0$$

$$\sum_j \frac{\pi_{ij}^i}{\mu_i} = \frac{\mu_i}{\mu_i} = 1$$

$$\sum_i \frac{\pi_{ij}^i}{\mu_i} \mu_i = v_j$$

P stochastic matrix ( $P \geq 0, P^T 1 = 1$ )

Frobenius-Perron:  $P$  pos. diagonals, irrecl.  
 $\rightarrow P$  has eigenvector  $e$  to largest (simple) eigenvector 1  
 $\lim_{n \rightarrow \infty} P^n x = e, \|x\|_1 = \|e\|_1$

**Disintegration theorem** in OT (book Ambrosio et al: Thm 5.3.1): existence of a Markov kernel  $K$  (uniquely defined for a.e.  $x \in X$ ) such that  $\pi = \mu \times K$  Usual notation  $K(x, \cdot) = \pi_x$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Remark on Transfer Operators

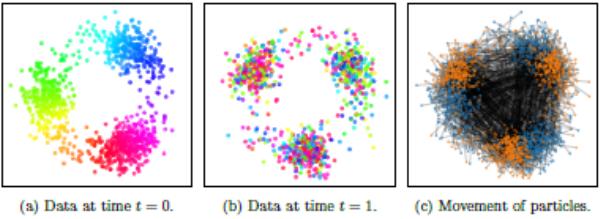


Figure 7: Particles moving in a potential with circular driving force. The color scheme illustrates the particle mixing.

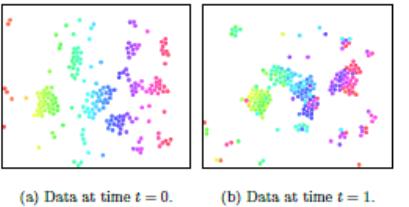
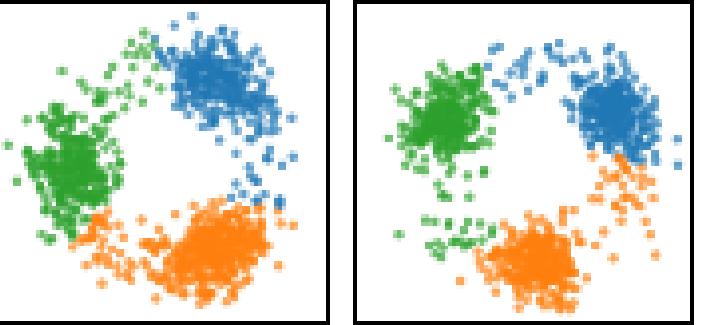
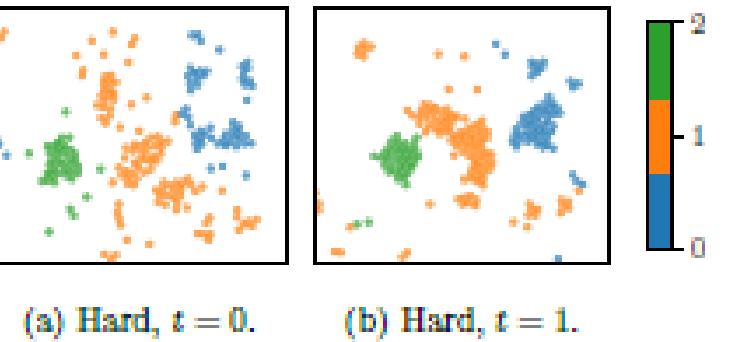


Figure 10: Particles moving in a potential with circular driving force. Again, the color scheme illustrates the particle mixing.



- Refs:
1. P Koltai, J von Lindheim, S Neumayer, G S: Transfer operators from optimal transport plans for coherent set detection Physica D, 2021
  2. F. Beier: Gromov–Wasserstein Transfer Operators
  3. Junge, O., Matthes, D., Schmitzer, B.: Entropic transfer operators. 2022 (exact transfer operator is known on a finite subset of the full state space. Then, using regularized OT, a finite-dimensional approximation is constructed which limit is a regularized version of the ground truth and exhibits desirable properties, such as retention of the spectral information)

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

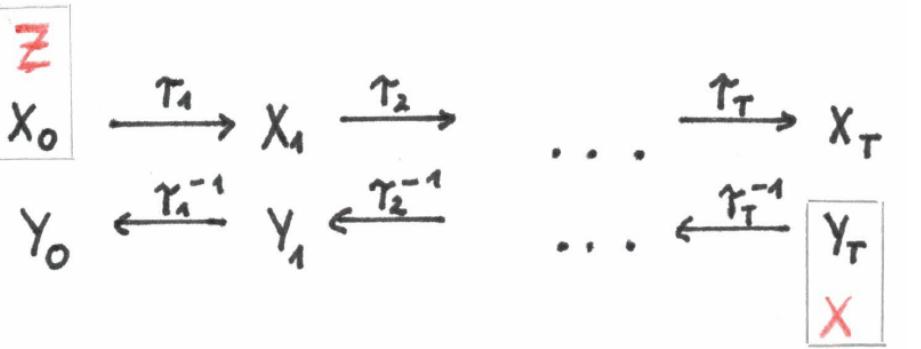
# Outline

1. Markov Kernels and Markov Chains
2. Normalizing Flows via Markov Chains
3. Generalized Normalizing Flows (GNFs)
4. Conditional GNFs

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Normalizing Flows via Markov Chains

Normalizing Flow:  $\mathcal{T}(\cdot; \theta) = \mathcal{T}_T \circ \dots \circ \mathcal{T}_1$ ,



- ◆  $\mathcal{T}_t: \mathbb{R}^d \rightarrow \mathbb{R}^d$  generate a pair of Markov chains  $((\textcolor{red}{X}_0, \dots, X_T), (\textcolor{blue}{Y}_T, \dots, Y_0))$

$$X_0 \sim P_Z, \quad X_t = \mathcal{T}_t(X_{t-1}) \quad \text{and}$$

$$Y_T \sim P_X, \quad Y_{t-1} = \mathcal{T}_t^{-1}(Y_t).$$

- ◆ Markov kernels  $\mathcal{K}_t(x, \cdot) = P_{X_t|X_{t-1}=x} = \delta_{\mathcal{T}_t(x)}$  and  $\mathcal{R}_t(x, \cdot) = \delta_{\mathcal{T}_t^{-1}(x)}$
- ◆ Minimize the „whole path”

$$\begin{aligned} \mathcal{L}_{\text{NF}}(\theta) &= \text{KL}(P_X, P_{\mathcal{T}_\# Z}) = \text{KL}(P_{Y_T}, P_{X_T}) \\ &= \text{KL}(P_{(Y_0, \dots, Y_T)}, P_{(X_0, \dots, X_T)}) \quad \text{in general } \leq \end{aligned}$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Normalizing Flows via Markov Chains

Deterministic Markov kernels:

$$\mathcal{K}_t(x, \cdot) = P_{X_t|X_{t-1}=x} = \delta_{\mathcal{T}_t(x)}$$

$$\mathcal{R}_t(y, \cdot) = P_{Y_{t-1}|Y_t=y} = \delta_{\mathcal{T}_t^{-1}(y)}$$

## Proof

$$\begin{aligned} P_{(X_{t-1}, X_t)}(A \times B) &= \int_{\mathbb{R}^d \times \mathbb{R}^d} 1_{A \times B}(x_{t-1}, x_t) dP_{(X_{t-1}, X_t)}(x_{t-1}, x_t) \\ &= \int_{\mathbb{R}^d \times \mathbb{R}^d} 1_A(x_{t-1}) 1_B(x_t) dP_{(X_{t-1}, X_t)}(x_{t-1}, x_t). \end{aligned}$$

Since  $P_{(X_{t-1}, X_t)}$  is by definition concentrated on the set  $\{(x_{t-1}, \mathcal{T}_t(x_{t-1})) : x_{t-1} \in \mathbb{R}^d\}$ , this becomes

$$\begin{aligned} P_{(X_{t-1}, X_t)}(A \times B) &= \int_{\mathbb{R}^d \times \mathbb{R}^d} 1_A(x_{t-1}) 1_B(\mathcal{T}_t(x_{t-1})) dP_{(X_{t-1}, X_t)}(x_{t-1}, x_t) \\ &= \int_A 1_B(\mathcal{T}_t(x_{t-1})) dP_{X_{t-1}}(x_{t-1}) \\ &= \int_A \delta_{\mathcal{T}_t(x_{t-1})}(B) dP_{X_{t-1}}. \end{aligned}$$

Consequently, the transition kernel  $\mathcal{K}_t = P_{X_t|X_{t-1}}$  is given by  $\mathcal{K}_t(x, \cdot) = \delta_{\mathcal{T}_t(x)}$ . □

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Outline

1. Markov Kernels and Markov Chains
2. Normalizing Flows via Markov Chains
3. Generalized Normalizing Flows (GNFs)
4. Conditional GNFs

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Generalized Normalizing Flows

**Generalized/Stochastic normalizing flow (SNF) = pair of Markov chains**

$$((X_0, \dots, X_T), (Y_T, \dots, Y_0))$$

that minimizes the loss function

$$\mathcal{L}_{\text{SNF}}(\theta) = \text{KL}(P_{(Y_0, \dots, Y_T)}, P_{(X_0, \dots, X_T)}) :$$

P1)  $P_{X_t}, P_{Y_t}$  have the densities  $p_{X_t}, p_{Y_t} : \mathbb{R}^{d_t} \rightarrow \mathbb{R}_{>0}$  for any  $t = 0, \dots, T$ .

P2) There exist Markov kernels  $\mathcal{K}_t = P_{X_t|X_{t-1}}$  and  $\mathcal{R}_t = P_{Y_{t-1}|Y_t}$ ,  $t = 1, \dots, T$ :

$$P_{(X_0, \dots, X_T)} = P_{X_0} \times P_{X_1|X_0} \times \cdots \times P_{X_T|X_{T-1}},$$

$$P_{(Y_T, \dots, Y_0)} = P_{Y_T} \times P_{Y_{T-1}|Y_T} \times \cdots \times P_{Y_0|Y_1}.$$

P3) For  $P_{X_t}$ -almost every  $x \in \mathbb{R}^{d_t}$ , the measures  $P_{Y_{t-1}|Y_t=x}$  and  $P_{X_{t-1}|X_t=x}$  are **absolutely continuous with respect to each other**

**Important:** Conditional distributions  $P_{X_t|X_{t-1}}$  themselves **must not be absolutely continuous!**

Indeed not the case for NFs and MCMC layers.

**Avoid: UMOs !**    e.g.  $\frac{\delta(x_t - T_t(x_{t-1}))}{\delta(x_{t-1} - T_t^{-1}(x_t))}$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

**Loss function:**  $\mathcal{L}_{\text{SNF}}(\theta) = \text{KL}(P_{(Y_0, \dots, Y_T)}, P_{(X_0, \dots, X_T)})$

**Theorem:** With Radon-Nikodym derivative  $f_t(\cdot, x_t) = \frac{dP_{Y_{t-1}|Y_t=x_t}}{dP_{X_{t-1}|X_t=x_t}}$  we have

$$\begin{aligned}\mathcal{L}_{\text{SNF}}(\theta) &= \mathbb{E}_{(x_0, \dots, x_T) \sim P_{(Y_0, \dots, Y_T)}} \left[ \log \left( \frac{p_X(x_T)}{p_Z(x_0)} \prod_{t=1}^T \frac{f_t(x_{t-1}, x_t) p_{X_{t-1}(x_{t-1})}}{p_{X_t}(x_t)} \right) \right] \\ &= \mathbb{E}_{(x_0, \dots, x_T) \sim P_{(Y_0, \dots, Y_T)}} \left[ \log \left( \frac{p_X(x_T)}{p_{X_T}(x_T)} \prod_{t=1}^T f_t(x_{t-1}, x_t) \right) \right]\end{aligned}$$

The right-hand side can be computed for the different layers:

- ◆ NF layer:  $\frac{p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} = \frac{1}{|\nabla \mathcal{T}_t^{-1}(x_t)|}$  and  $f_t(x_{t-1}, x_t) = 1$
- ◆ MCMC layer:  $\frac{f_t(x_{t-1}, x_t) p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} = \frac{p_t(x_{t-1})}{p_t(x_t)}$
- ◆ Langevin layer:  $\frac{f_t(x_{t-1}, x_t) p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} = \exp \left( \frac{1}{2} (\|\eta_t\|^2 - \|\tilde{\eta}_t\|^2) \right)$ ,  
with proposal density  $p_t: \mathbb{R}^d \rightarrow \mathbb{R}_{>0}$ ,  $u_t = -\log(p_t)$  and  
 $\eta_t := \frac{1}{a_2} (x_{t-1} - x_t - a_1 \nabla u_t(x_{t-1}))$ ,  $\tilde{\eta}_t := \frac{1}{a_2} (x_{t-1} - x_t + a_1 \nabla u_t(x_t))$
- ◆ Diffusion layer:  $\frac{f_t(x_{t-1}, x_t) p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} = \exp \left( \frac{1}{2} (\|\eta_t\|^2 - \|\tilde{\eta}_t\|^2) \right)$ , where  
 $\eta_t := \frac{1}{\sqrt{\epsilon} h_{t-1}} (x_{t-1} - x_t + \epsilon g_{t-1}(x_{t-1}))$ ,  $\tilde{\eta}_t := \frac{1}{\sqrt{\epsilon} g_t} (x_{t-1} - x_t - \epsilon (\mathbf{g}_t(x_t) - \mathbf{h}_t^2 s_t(x_t)))$
- ◆ VAEs:  $\frac{f_t(x_{t-1}, x_t) p_{X_{t-1}}(x_{t-1})}{p_{X_t}(x_t)} = \frac{q_\phi(x_{t-1}|x_t)}{p_\theta(x_t|x_{t-1})}$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Learning SNFs

$$\int_{A \times B} f(x, y) dP_{(X, Y)} = \int_A \int_B f(x, y) dP_Y(y) dP_{X|Y=y}(x)$$

$$\textcircled{1} \quad \rho := \frac{dP_{(Y_1, \dots, Y_T)}}{dP_{(X_0, \dots, X_T)}} = \frac{\rho_{Y_T}(x_T)}{\rho_{X_T}(x_T)} \overbrace{\prod_{t=1}^T f_t(x_{t-1}, x_t)}^{\frac{dP_{Y_{t-1}|Y_t=x_t}}{dP_{X_{t-1}|X_t=x_t}}}$$

$$\textcircled{2} \quad KL(\mu, \nu) = \int \log \underbrace{\frac{\partial \mu}{\partial \nu}}_{\text{Radon-Nikodym derivative}} d\mu(x) = \mathbb{E}_{x \sim \mu} [\log \frac{\partial \mu}{\partial \nu}]$$

$$KL(P_{(Y_0, \dots, Y_T)}, P_{(X_0, \dots, X_T)}) = \mathbb{E}_{(X_0, \dots, X_T) \sim P_{(Y_0, \dots, Y_T)}} [\log \frac{\partial P_{(Y_0, \dots, Y_T)}}{\partial P_{(X_0, \dots, X_T)}}]$$

- \textcircled{1} can be shown by induction using •

Case  $T=1$  ;  $\rho$  r.h.s in \textcircled{1}

$$\begin{aligned} & \text{To show : } \rho \text{ of } P_{(X_0, X_1)} = \frac{dP_{(Y_0, Y_1)}}{dP_{(X_0, X_1)}} \\ & \int_{A \times B} \frac{P_{Y_1}(x_1)}{P_{X_1}(x_1)} \frac{dP_{Y_0|Y_1=x_1}}{dP_{X_0|X_1=x_1}} dP_{(X_0, X_1)}(x_0, x_1) \\ & = \int_A \int_B \underbrace{- - -}_{= \int_A \int_B dP_{Y_1}(x_1) dP_{Y_0|Y_1=x_1}(x_0)} \underbrace{dP_{X_1}(x_1) dP_{X_0|X_1=x_1}(x_0)}_{= \int_{A \times B} dP_{(Y_0, Y_1)}(x_0, x_1)} \\ & = \int_{A \times B} dP_{(Y_0, Y_1)}(x_0, x_1) \end{aligned}$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

- ◆ Markov chain:

$$X_t = X_{t-1} + 1_{[U,1]}(\alpha_t(X_{t-1}, X_{t-1} + \xi_t)) \xi_t, \quad \xi_t \sim \mathcal{N}(0, \sigma^2 I)$$

where

$$\alpha_t(x, y) := \min \left\{ 1, \frac{p_t(y)}{p_t(x)} \right\}$$

with proposal density  $p_t: \mathbb{R}^d \rightarrow \mathbb{R}_{>0}$

- ◆ Markov kernels  $\mathcal{K}_t = \mathcal{R}_t: \mathbb{R}^d \times \mathcal{B}(\mathbb{R}^d) \rightarrow [0, 1]$

$$\begin{aligned} \mathcal{K}_t(x, A) &:= \int_A \mathcal{N}(y; x, \sigma^2 I) \alpha_t(x, y) dy \\ &\quad + \delta_x(A) \int_{\mathbb{R}^d} \mathcal{N}(y; x, \sigma^2 I) (1 - \alpha_t(x, y)) dy \end{aligned}$$

- ◆ one sampling step of Metropolis-Hastings algorithm  
(known to sample from the proposal distribution  $p_t$ )

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Metropolis-Hastings Algorithm

**Input:**  $x_0 \in \mathbb{R}^d$ , proposal density  $p_t: \mathbb{R}^d \rightarrow \mathbb{R}_{>0}$

**For**  $k = 0, 1, \dots$  do

Draw  $x'$  from  $\mathcal{N}(x_k, \sigma^2 I)$  and  $u$  uniformly in  $[0, 1]$ .

Compute the acceptance ratio

$$\alpha(x_k, x') := \min \left\{ 1, \frac{p(x')}{p(x_k)} \right\}.$$

Set

$$x_{k+1} := \begin{cases} x' & \text{if } u < \alpha(x_k, x'), \\ x_k & \text{otherwise.} \end{cases}$$

**Output:**  $\{x_k\}_k$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# VAE Layer

1. Autoencoders (AE):  $\mathcal{L}_{\text{AE}}(\phi, \theta) := \mathbb{E}_{x \sim P_X} (\|x - D_\theta(E_\phi(x))\|^2)$

- encoder  $E = E_\phi: \mathbb{R}^d \rightarrow \mathbb{R}^n, d > n$
- decoder  $D = D_\theta: \mathbb{R}^n \rightarrow \mathbb{R}^d$

**Example:** Principle Component Analysis = affine encoder/decoder

Given  $x_i \in \mathbb{R}^m, i = 1, \dots, N$  (e.g. realizations of a random variable)

Find  $A \in \text{St}(m, d) \subset R^{m,d}, d \ll m$  i.e.  $A^\top A = I$  and  $b \in \mathbb{R}^m$  minimizing

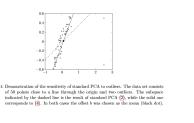
$$\min_{A \in \text{St}(m, d), b, t_i} \sum_{i=1}^N \|At_i + b - x_i\|^p, \quad p \in [1, \infty)$$

◆ Independent of  $p$ :  $t_i = (A^\top A)^{-1} A^\top (x_i - b) = A^\top (x_i - b), i = 1, \dots, N$  Thus

$$\min_{A \in \text{St}(m, d), b} \sum_{i=1}^N \|AA^\top(x_i - b) + b - x_i\|^p = \min_{A \in \text{St}(m, d), b} \sum_{i=1}^N \|(I - AA^\top)(b - x_i)\|^p$$

◆  $p = 1$ : robust PCA (in particular computation of  $b$  not trivial ??)

Neumayer, Nimmer, Setzer, Steidl: On the robust PCA and Weiszfeld's algorithm + On the rotational invariant L1-norm PCA (2020)



1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Autoencoders

$p = 2$ :

$$\min_{A \in \text{St}(m,d), b} \sum_{i=1}^N \| \underbrace{(I - AA^\top)}_{P \text{ orth. proj.}} (b - x_i) \|^2$$

- ◆ Computation of  $b$ :  $\sum_{i=1}^N P(b - x_i) = 0 \Rightarrow b = \sum_{i=1}^N \frac{1}{N} x_i + \mathcal{R}(A) = \bar{x} + \mathcal{R}(A)$
- ◆ Computation of  $A$ :  $y_i = x_i - \bar{x}$

$$\min_{A \in \text{St}(m,d)} \sum_{i=1}^N \|y_i - AA^\top y_i\|^2 = \min_{A \in \text{St}(m,d)} \sum_{i=1}^N \|y_i\|^2 - \langle AA^\top y_i, y_i \rangle$$

Find

$$\begin{aligned} & \operatorname{argmax}_{A \in \text{St}(m,d)} \sum_{i=1}^N y_i^\top A A^\top y_i \\ &= \operatorname{argmax}_{A \in \text{St}(m,d)} \sum_{k=d}^N a_k^\top \underbrace{Y Y^\top}_{C = \text{Cov}(X)} a_k \end{aligned}$$

Find  $d$  orthogonal eigenvectors  $a_1, \dots, a_d$  of  $C$  belonging to the largest  $d$  eigenvectors of  $C$

Autoencoder:

$$\begin{aligned} \|x - D_\theta(E_\phi(x))\|^2 &= \|x - b - AA^\top(x - b)\|^2 \\ E_\phi(x) &:= A^\top(x - b), \quad D_\theta(y) = Ay + b \end{aligned}$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Variational Autoencoders

## 2. Variational Autoencoders (VAE):

- encoder  $D(z) = D_\theta(z) := (\mu_\theta(z), \Sigma_\theta(z))$

$$P_{X_1|Z=z}(\cdot) = \mathcal{K}(z, \cdot) := \mathcal{N}(\mu_\theta(z), \Sigma_\theta(z)), \quad p_\theta(x|z) = \mathcal{N}(x; \mu_\theta(z), \Sigma_\theta(z))$$

- decoder  $E(x) = E_\phi(x) := (\mu_\phi(x), \Sigma_\phi(x))$

$$P_{Y_0|X=x}(\cdot) = \mathcal{R}(x, \cdot) := \mathcal{N}(\mu_\phi(x), \Sigma_\phi(x)), \quad q_\phi(z|x) = \mathcal{N}(z; \mu_\phi(x), \Sigma_\phi(x))$$

SNF  $((X_0, X_1), (Y_1, Y_0))$  one layer Markov chain with

$$\begin{aligned} \mathcal{L}_{\text{SNF}}(\theta, \phi) &= \mathbb{E}_{(z, x) \sim P_{(Y_0, Y_1)}} \left[ \log \left( \frac{p_X(x) f_1(z, x)}{p_{X_1}(x)} \right) \right] \\ &= - \underbrace{\mathbb{E}_{x \sim P_X} [\mathcal{L}_{\theta, \phi}(x)]}_{\text{ELBO}} + \underbrace{\mathbb{E}_{x \sim P_X} [\log(p_X(x))]}_{\text{const}} \end{aligned}$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

Markov Chain (2 layers):

$$P_{Y_0|X_1=X}(z) = \frac{P(Y_0|X_1=X)}{P(Z|X_1=X)} = \frac{q_\phi(z|x)}{p_\theta(x_1=x)} = \frac{q_\phi(z|x)}{p_\theta(x_1=x)} = \frac{q_\phi(z|x)}{p_\theta(x_1=x)} = \frac{q_\phi(z|x)}{p_\theta(x_1=x)}$$

Firm to learn  $\theta$  such that

$$P_X(x) \approx \int_R p_\theta(x|z) p(z) dz$$

$$P_X(x) \approx \int_R q_\phi(x|z) p(z) dz$$

- Sample from  $P_X$   $\xrightarrow{\text{evidence}}$  sample from  $p_\theta(x)$
- Loss function: Idea:  $E_{X \sim P_X} [\log \frac{p_\theta(x)}{q_\phi(x)}]$   $\rightarrow$  max

Approximation of evidence from below: evident lower bound (ELBO)

$$\begin{aligned} & \log(p_\theta(x)) = E_{Z \sim q_\phi(\cdot|x)} [\log \frac{p_\theta(x|z)}{p_\theta(x)}] \\ & = E_{Z \sim q_\phi(\cdot|x)} \left[ \log \frac{q_\phi(x|z)}{q_\phi(x)} \right] + E_{Z \sim q_\phi(\cdot|x)} [\log \frac{q_\phi(x|z)}{p_\theta(x|z)}] \\ & \xrightarrow{\text{by def.}} E_{Z \sim q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x|z)}{q_\phi(x|z)} \right] + KL(p_\theta(x|z), q_\phi(x|z)) \\ & \geq E_{Z \sim q_\phi(\cdot|x)} [\log \frac{p_\theta(x|z)}{q_\phi(x|z)}] \end{aligned}$$

Relation to Markov chains:  $L_{\text{SNT}}(\theta, \psi) = KL(P_{(Y_0, Y_1)}, P_{(X_0, X_1)})$

$$\begin{aligned} & = E_{(X_0, X_1) \sim P_{(Y_0, Y_1)}} \left[ \log \frac{p_{Y_0, Y_1}(x_0, x_1)}{p_{X_0, X_1}(x_0, x_1)} \right] \\ L_{\text{SNT}}(\theta, \psi) & = E_{(Z, X) \sim P_{(Y_0, Y_1)}} \left[ \log \frac{p_{Y_0, Y_1}(x)}{p_{X, Y}(x)} \right] \end{aligned}$$

Thus

$$\begin{aligned} L_{\text{SNT}}(\theta, \psi) &= E_{(Z, X) \sim P_{(Y_0, Y_1)}} \left[ \log \frac{q_\phi(z|x)}{p_\theta(x|z)} \right] \\ P_{(Y_0, Y_1)} &= P_X \times \underbrace{P_{Y_0|X}}_{p_\theta(x)} \quad (\text{Markov kernel property}) \\ &= \mathbb{E}_{X \sim P_X} \left( \mathbb{E}_{Z \sim q_\phi(\cdot|x)} \left[ \log \frac{q_\phi(z|x)}{p_\theta(x|z)} \right] \right) \\ &= \mathbb{E}_{X \sim P_X} \left( \mathbb{E}_{Z \sim q_\phi(\cdot|x)} \left[ \log \frac{q_\phi(z|x)}{p_\theta(x|z)} \right] - ELBO \right. \\ &\quad \left. - \mathbb{E}_{X \sim P_X} [\log p_\theta(x)] \right) \\ &\quad \xrightarrow{\text{const.}} \\ &\Rightarrow \underset{\theta, \psi}{\text{minimize}} \ L_{\text{SNT}}(\theta, \psi) = \underset{\theta, \psi}{\text{argmax}} \ E_{X \sim P_X} (ELBO) \end{aligned}$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Langevin Layer

Markov chain:

$$X_t := X_{t-1} + a_1 \nabla \log p_t(X_{t-1}) + a_2 \xi_t,$$

where  $a_1, a_2 > 0$  constants,  $\xi_t \sim \mathcal{N}(0, I)$  and  $p_t$  proposal density

Markov kernels  $\mathcal{K}_t = \mathcal{R}_t: \mathbb{R}^d \times \mathcal{B}(\mathbb{R}^d) \rightarrow [0, 1]$ :

$$\mathcal{K}_t(x, \cdot) = \mathcal{N}(x - a_1 \nabla u_t(x), a_2^2 I).$$

**Proof** Use independence of  $\xi_t$  of  $X_t$  and  $X_{t-1}$  to obtain that  $X_t$  and  $X_{t-1}$  have the common density

$$\begin{aligned} p_{(X_{t-1}, X_t)}(x_{t-1}, x_t) &= p_{X_{t-1}, \xi_t}(x_{t-1}, \frac{1}{a_2}(x_t - x_{t-1} + a_1 \nabla u_t(x_{t-1}))) \\ &= p_{X_{t-1}}(x_{t-1}) p_{\xi_t}(\frac{1}{a_2}(x_t - x_{t-1} + a_1 \nabla u_t(x_{t-1}))) \\ &= p_{X_{t-1}}(x_{t-1}) \mathcal{N}(x_t; x_{t-1} - a_1 \nabla u_t(x_{t-1}), a_2^2 I). \end{aligned}$$

Then, for  $A, B \in \mathcal{B}(\mathbb{R}^d)$ , it holds

$$\begin{aligned} P_{(X_{t-1}, X_t)}(A \times B) &= \int_{A \times B} p_{X_{t-1}}(x_{t-1}) \mathcal{N}(x_t; x_{t-1} - a_1 \nabla u_t(x_{t-1}), a_2^2 I) d(x_{t-1}, x_t) \\ &= \int_A \int_B \mathcal{N}(x_t; x_{t-1} - a_1 \nabla u_t(x_{t-1}), a_2^2 I) dx_t p_{X_{t-1}}(x_{t-1}) dx_{t-1} \\ &= \int_A \mathcal{K}_t(x_{t-1}, B) p_{X_t}(x_t) dP_{X_{t-1}}(x_{t-1}). \end{aligned}$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Background: Langevin Layer (for fixed $t$ )

Overdamped Langevin SDE

$$X(0) = X^0,$$

$$dX(\tau) = -\nabla \Psi(X(\tau))d\tau + \sqrt{2\beta^{-1}}dW(\tau)$$

Euler-Maruyama forward step:

$$\begin{aligned} X_t &:= X_{t-1} - \eta \nabla \Psi(X_{t-1}) + \sqrt{2\beta^{-1}\eta} \xi_t \\ &= X_{t-1} + \underbrace{\eta \beta^{-1}}_{a_1} \nabla \log p(X_{t-1}) + \underbrace{\sqrt{2\beta^{-1}\eta}}_{a_2} \xi_t \end{aligned}$$

where  $p(x) = C^{-1}e^{-\beta\Psi(x)}$ ,  $C$  normalizing factor

◆  $\nabla \log p$  is known as **score**

Density of random variables is described by Fokker-Planck equation

= Wasserstein gradient flow of  $\mathcal{F}(\mu) = \text{KL}(\mu, p dx)$

$$\rho(x, 0) = \rho^0(x)$$

$$\frac{\partial \rho}{\partial \tau} = \text{div}(\rho \nabla \Psi(x)) + \beta^{-1} \Delta \rho$$

As  $\tau \rightarrow \infty$  stationary solution:  $p(x) = C^{-1}e^{-\beta\Psi(x)}$ ,  $C$  normalizing factor

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

## Relation to Proximal Operator

For a proper, lsc function  $\mathcal{F}: \mathbb{P}_2(\mathbb{R}^d) \rightarrow (-\infty, \infty]$  and  $\tau > 0$ , the **Wasserstein proximal mapping** is the set-valued function

$$\text{prox}_{\tau\mathcal{F}}(\mu) := \operatorname{argmin}_{\nu \in \mathbb{P}_2(\mathbb{R}^d)} \left\{ \frac{1}{2\tau} W_2^2(\mu, \nu) + \mathcal{F}(\nu) \right\}, \quad \mu \in \mathbb{P}_2(\mathbb{R}^d).$$

- ◆ existence and uniqueness of the minimizer is assured if  $\mathcal{F}$  is  $\lambda$ -convex along generalized geodesics, where  $\lambda > -1/\tau$  and  $\mu \in \text{dom } \mathcal{F}$

**Backward scheme = Jordan-Kinderlehrer-Otto (JKO) scheme** starting at  $\mu_\tau^0 := \mu \in \mathbb{P}_2(\mathbb{R}^d)$  with time step size  $\tau$  is the curve  $\gamma_\tau: [0, +\infty) \rightarrow \mathbb{P}_2(\mathbb{R}^d)$  given by

$$\gamma_\tau|_{((n-1)\tau, n\tau]} := \mu_\tau^n := \text{prox}_{\tau\mathcal{F}}(\mu_\tau^{n-1})$$

- ◆ If  $\mathcal{F}: \mathbb{P}_2(\mathbb{R}^d) \rightarrow (-\infty, +\infty]$  is coercive and  $\lambda$ -convex along generalized geodesics, then the JKO curves  $\gamma_\tau$  starting at  $\mu \in \overline{\text{dom } \mathcal{F}}$  converge for  $\tau \rightarrow 0$  locally uniformly to a locally Lipschitz curve  $\gamma: (0, +\infty) \rightarrow \mathbb{P}_2(\mathbb{R}^d)$  which is the unique Wasserstein gradient flow of  $\mathcal{F}$  with  $\gamma(0+) = \mu$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Diffusion Layer

**SDE:**  $dX_t = g_t(X_t)dt + h_t dW_t$

Explicit Euler discretization with step size  $\epsilon > 0$ :

$$X_t = X_{t-1} + \epsilon g_{t-1}(X_{t-1}) + \sqrt{\epsilon} h_{t-1} \xi_{t-1}, \quad t = 1, \dots, T,$$

where  $\xi_{t-1} \sim \mathcal{N}(0, I)$  is independent of  $X_0, \dots, X_{t-1}$

Markov kernel:

$$\mathcal{K}_t(x, \cdot) = P_{X_t|X_{t-1}=x} = \mathcal{N}(x + \epsilon g_{t-1}(x), \epsilon h_{t-1}^2).$$

- ◆ Song et al. 2020 parametrized the functions  $g_t: \mathbb{R}^d \rightarrow \mathbb{R}^d$  by some a-priori learned **scoring network**
- ◆ Motivated by the time-reversal process of the SDE, Zhang/Chen 2021 introduced the **backward layer**

$$\mathcal{R}_t(x, \cdot) = P_{Y_{t-1}|Y_t=x} = \mathcal{N}(x + \epsilon(g_t(x) - h_t^2 s_t(x)), \epsilon h_t^2)$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Tweedie Formula

$$Y = x + \eta, \eta \sim \mathcal{N}(0, \sigma^2 I)$$

$$p(y) = \int p(y|x)p(x) dx = C \int e^{-\|y-x\|^2/(2\sigma^2)} p(x) dx$$

MMSE (maximum mean square error):

$$\begin{aligned} \hat{x}(y) &= \mathbb{E}[X|Y=y] = \int xp(x|y) dx = \int x \frac{p(y|x)p(x)}{p(y)} dx \\ &= y + \sigma^2 \nabla_y \log p(y) \quad \text{Tweedie formula} \end{aligned}$$

since

$$\begin{aligned} \nabla_y p(y) &= C \int \frac{1}{\sigma^2} e^{-\|y-x\|^2/(2\sigma^2)} (x - y) p(x) dx \\ \frac{\sigma^2}{p(y)} \nabla_y p(y) &= \frac{C}{p(y)} \int e^{-\|y-x\|^2/(2\sigma^2)} x p(x) dx - \frac{Cy}{p(y)} \int e^{-\|y-x\|^2/(2\sigma^2)} p(x) dx \\ &= \underbrace{\int p(x|y)x dx}_{\hat{x}(y)} - y \underbrace{\int p(x|y) dx}_1 \\ \hat{x}(y) - y &= \frac{\sigma^2}{p(y)} \nabla_y p(y) = \sigma^2 \nabla_y \log p(y) \end{aligned}$$

Refs: papers of Kadkhodaie, Simoncelli et al. 2020,

Laumont, Almansa, Pereyra, Delon et al.: Bayesian imaging using PnP priors: when Langevin meets Tweedie, 2022

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Tweedie Formula

Special case used in MMSE denoising:

$$Y = X + \eta, \quad \eta \sim \mathcal{N}(0, \sigma^2 I), \quad X \sim \mathcal{N}(\mu_x, \Sigma_x) \quad \Rightarrow \quad \Sigma_Y = \Sigma_X + \sigma^2 I$$

Then

$$\hat{x}(y) = \mu_y + (\Sigma_y - \sigma^2 I) \Sigma_Y^{-1} (y - \mu_y)$$

Application:

- ◆ **MMSE estimation for similar patches:** Choose  $s \times s$  neighborhood (patch)  $y_i$  centered at  $i = (i_1, i_2) \in \mathcal{G}$  and interpret this and similar patches as realization of an  $p = s^2$ -dimensional random vector  $Y_i \sim \mathcal{N}(\mu_i, \Sigma_i)$

$$\hat{y}_j = \hat{\mu}_i + (\hat{\Sigma}_i - \sigma^2 I_p) \hat{\Sigma}_i^{-1} (y_j - \hat{\mu}_i), \quad j \in \mathcal{S}(i).$$

where  $\mathcal{S}(i)$  is the set of centers of patches similar to  $y_i$

- ◆ More fun on manifolds: **MMSE estimation for manifold-valued images:**

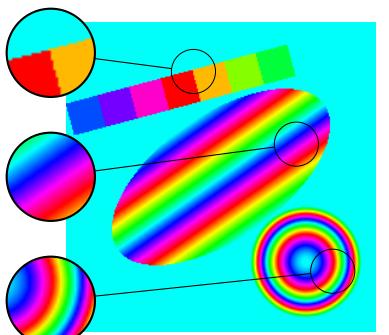
$$\hat{y}_j = \exp_{\hat{\mu}_i} \left( (\hat{\Sigma}_i - \sigma^2 I_{pd}) \hat{\Sigma}_i^{-1} (\log_{\hat{\mu}_i} y_j) \right), \quad j \in \mathcal{S}(i)$$

Refs: MMSE by Lebrun/Morel (Denoising cuisine) 2013,

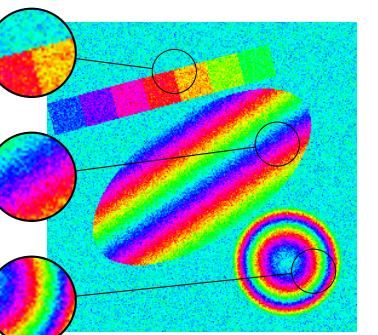
Laus, Nikolova, Persch, Steidl: A nonlocal denoising algorithm for manifold-valued images using second order statistics 2017

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

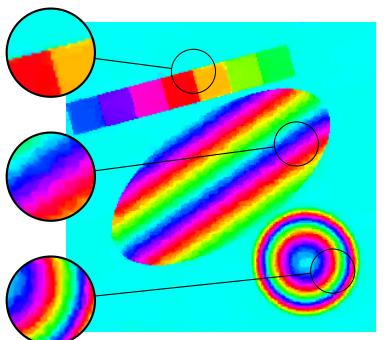
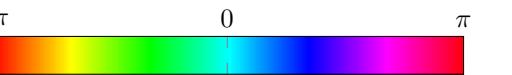
## Examples



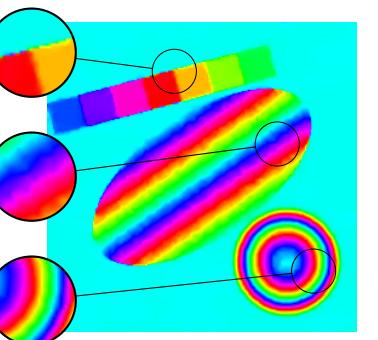
Original



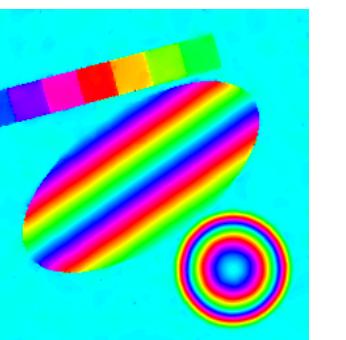
Noisy ( $88.5 \times 10^{-3}$ )



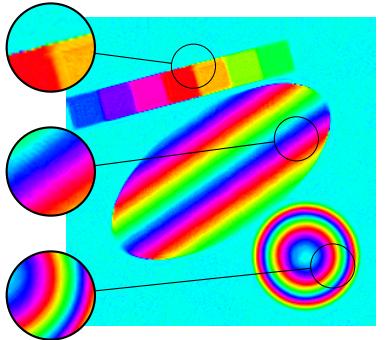
TV ( $7.2 \times 10^{-3}$ )



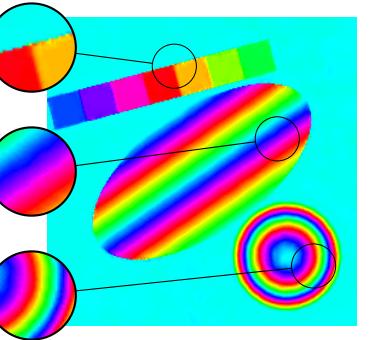
TV-TV2 ( $5.2 \times 10^{-3}$ )



TGV ( $2.6 \times 10^{-3}$ )



NL-means ( $8.1 \times 10^{-3}$ )



NL-MMSE ( $2.5 \times 10^{-3}$ )

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# Outline

1. Markov Kernels and Markov Chains
2. Normalizing Flows via Markov Chains
3. Generalized Normalizing Flows (GNFs)
4. Conditional GNFs

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

## 4. Conditional Stochastic Normalizing Flows

Inverse Problems:

$$Y = F(x) + \Xi, \quad \Xi \sim \mathcal{N}(0, \sigma) \Rightarrow Y \sim \mathcal{N}(F(x), \sigma)$$

$$Y = F(X) + \Xi$$

**Aim:** Sample from  $P_{X|Y=y}$  for some  $y$

**Idea:**  $P_{X|Y=y} \approx \mathcal{T}(y, \cdot)_\# P_Z$

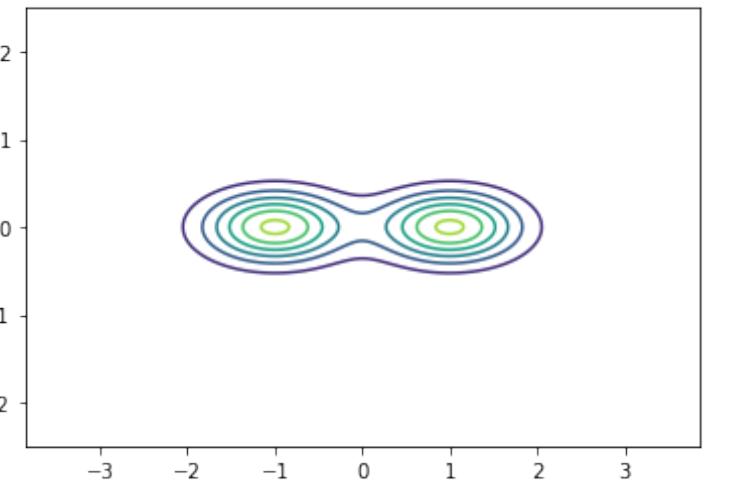
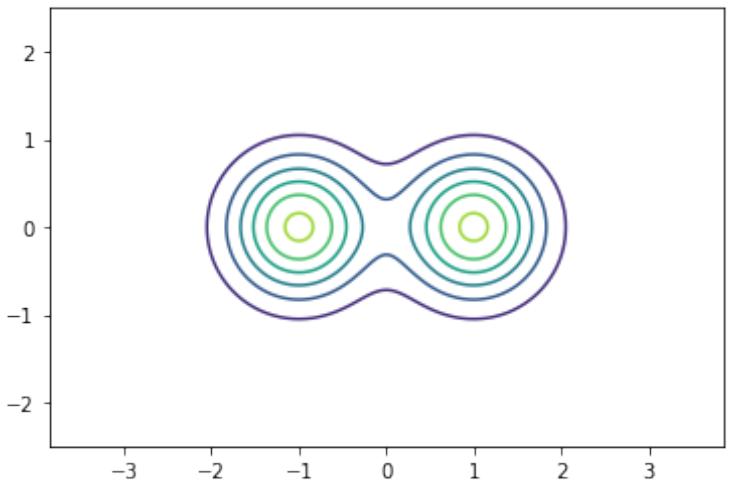


Illustration of the prior density  $p_X$  (left) and the posterior density  $p_{X|Y=y}$  for  $y = 0$  (right) within the inverse problem with  $F(x_1, x_2) = x_2$  and  $\Xi \sim \mathcal{N}(0, 0.1^2)$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

Conditional SNF conditioned to  $Y = y$  is pair of Markov chains

$$((X_0, \dots, X_T), (Y_T, \dots, Y_0))$$

cP1) the conditional distributions  $P_{X_t|Y=y}$  and  $P_{Y_t|Y=y}$  have densities

$$p_{X_t}(y, \cdot): \mathbb{R}^{d_t} \rightarrow \mathbb{R}_{>0}, \quad p_{Y_t}(y, \cdot): \mathbb{R}^{d_t} \rightarrow \mathbb{R}_{>0}$$

cP2) for  $P_Y$ -almost every  $y$ , there exist Markov kernels

$\mathcal{K}_t: \mathbb{R}^{\tilde{d}} \times \mathbb{R}^d \times \mathcal{B}(\mathbb{R}^d) \rightarrow [0, 1]$  and  $\mathcal{R}_t: \mathbb{R}^{\tilde{d}} \times \mathbb{R}^d \times \mathcal{B}(\mathbb{R}^d) \rightarrow [0, 1]$  such that

$$P_{(X_0, \dots, X_T)|Y=y} = P_{X_0} \times \mathcal{K}_1(y, \cdot, \cdot) \times \cdots \times \mathcal{K}_T(y, \cdot, \cdot),$$

$$P_{(Y_T, \dots, Y_0)|Y=y} = P_{Y_T} \times \mathcal{R}_T(y, \cdot, \cdot) \times \cdots \times \mathcal{R}_1(y, \cdot, \cdot).$$

cP3) for  $P_{Y, X_t}$ -almost every pair  $(y, x) \in \mathbb{R}^{\tilde{d}} \times \mathbb{R}^d$ , the measures  $P_{Y_{t-1}|Y_t=x, Y=y}$  and  $P_{X_{t-1}|X_t=x, Y=y}$  are absolute continuous with respect to each other.

**Loss Function:**  $\mathcal{L}_{\text{cSNF}}(\theta) = \text{KL}(P_{Y, (Y_0, \dots, Y_T)}, P_{Y, (X_0, \dots, X_T)})$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

# 1. Approximation of the Posterior for Gaussian Mixtures

- ◆ Posterior distribution is analytically known = ground truth

**Lemma:** Let  $X \sim \sum_{k=1}^K w_k \mathcal{N}(m_k, \Sigma_k)$ . Suppose that

$$Y = AX + \Xi, \quad A : \mathbb{R}^d \rightarrow \mathbb{R}^{\tilde{d}}, \quad \Xi \sim N(0, b^2 I)$$

Then

$$P_{X|Y=y} \propto \sum_{k=1}^K \tilde{w}_k \mathcal{N}(\cdot | \tilde{m}_k, \tilde{\Sigma}_k),$$

where

$$\tilde{\Sigma}_k := (\frac{1}{b^2} A^\top A + \Sigma_k^{-1})^{-1}, \quad \tilde{m}_k := \tilde{\Sigma}_k (\frac{1}{b^2} A^\top y + \Sigma_k^{-1} \mu_k).$$

and

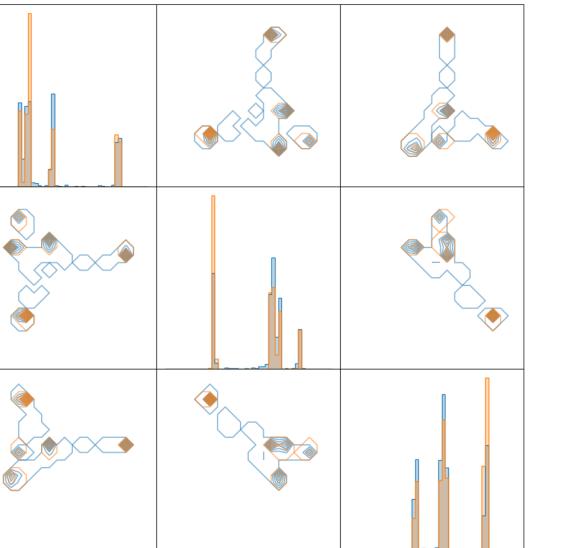
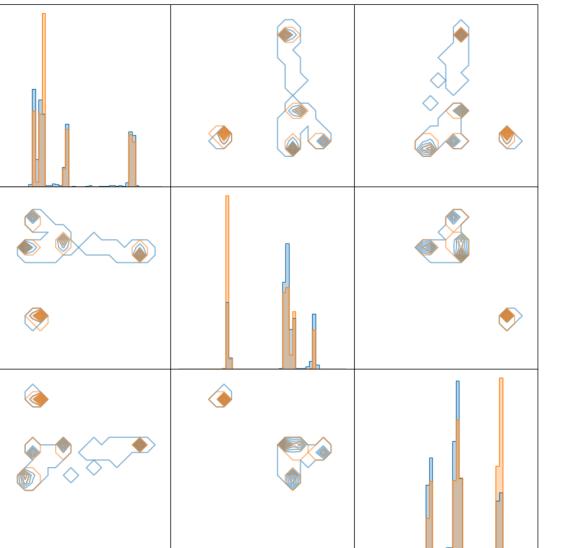
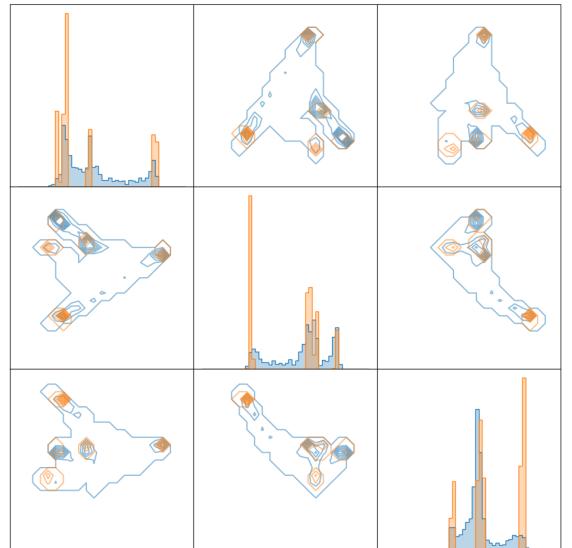
$$\tilde{w}_k := \frac{w_k}{|\Sigma_k|^{\frac{1}{2}}} \exp \left( \frac{1}{2} (\tilde{m}_k \tilde{\Sigma}_k^{-1} \tilde{m}_k - m_k \Sigma_k^{-1} m_k) \right).$$

**Experiment:** with GMM in  $\mathbb{R}^{100}$  with  $K = 5$  mixture components

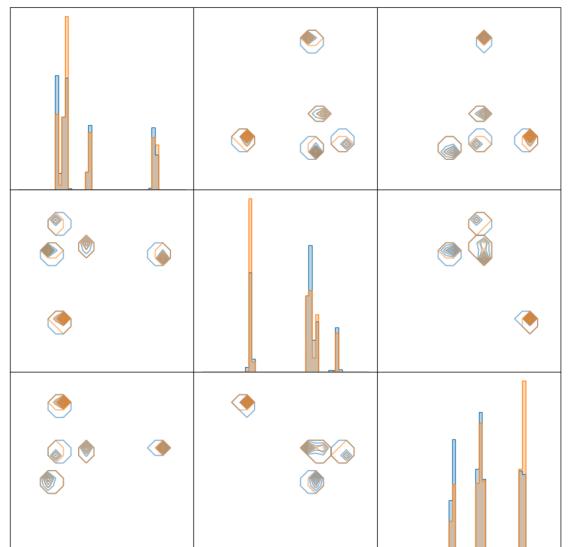
- ◆  $A$  diagonal matrix
- ◆  $\Xi \sim N(0, 0.05I)$
- ◆ proposal densities:  $p_t^y(x) = c_y (p_Z(x) p_{X|Y=y}(x))^{1/2}$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

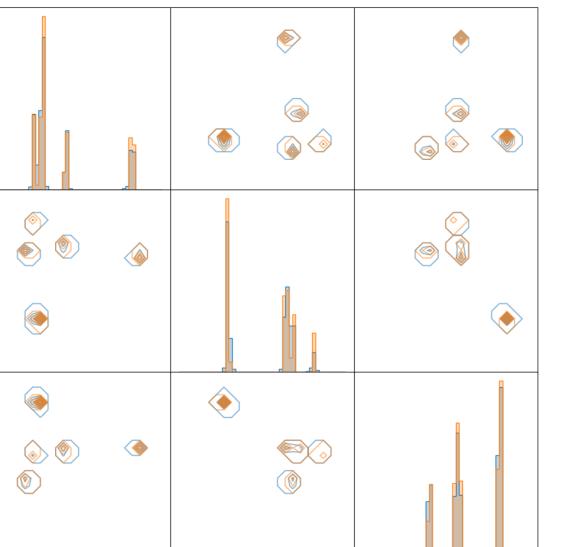
# 1. Approximation of the Posterior for Gaussian Mixtures



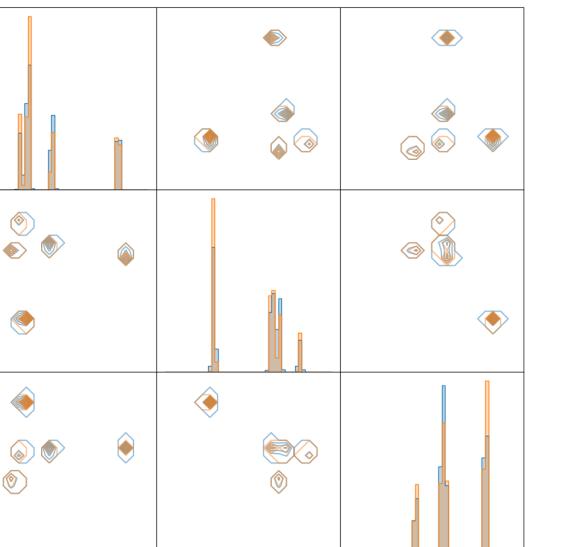
NF



NF + MALA



VAE



NF + VAE

Method	NF	NF+MALA	VAE	VAE+NF	VAE+MALA	NF+VAE+MALA
$W_1$	3.55	2.92	1.22	1.18	0.8	0.82

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

## 2. Example from Scatterometry

The parameters in  $x$ -space describe the geometry of the photo masks and

$$Y = F(X) + \eta$$

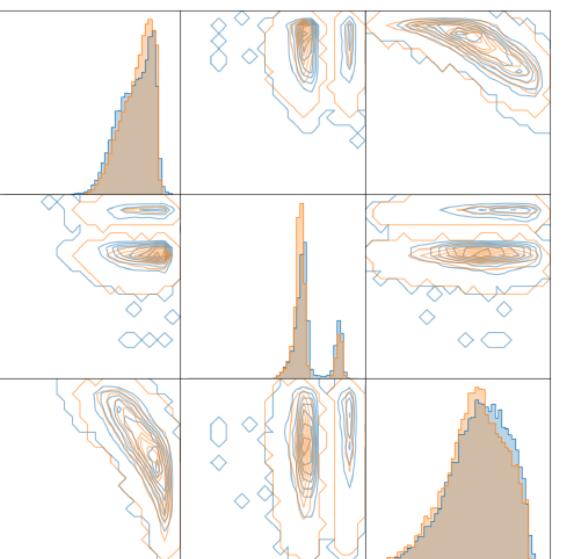
the observed diffraction pattern of light, where

- $F: \mathbb{R}^3 \rightarrow \mathbb{R}^{23}$  (from a nonlinear PDE and learned by NN)
- multiplicative + additive noise model

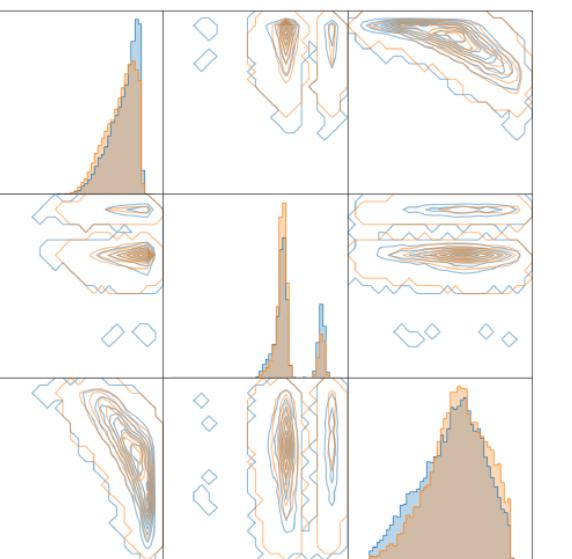
$$\eta = aF(X)\eta_1 + b\eta_2, \quad \eta_1, \eta_2 \sim \mathcal{N}(0, I), \quad a, b > 0$$

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

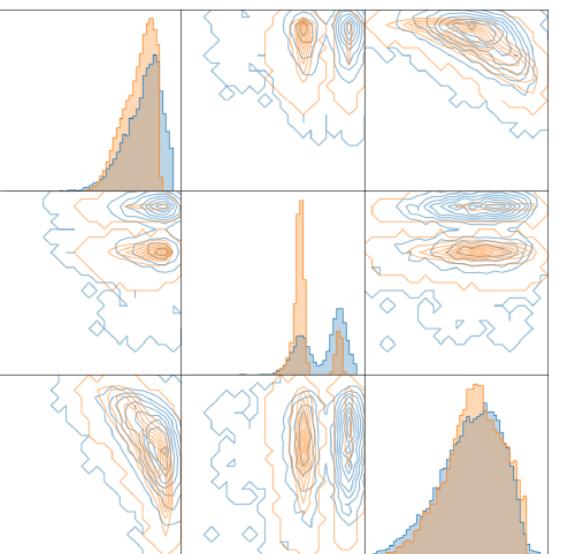
## 2. Example from Scatterometry



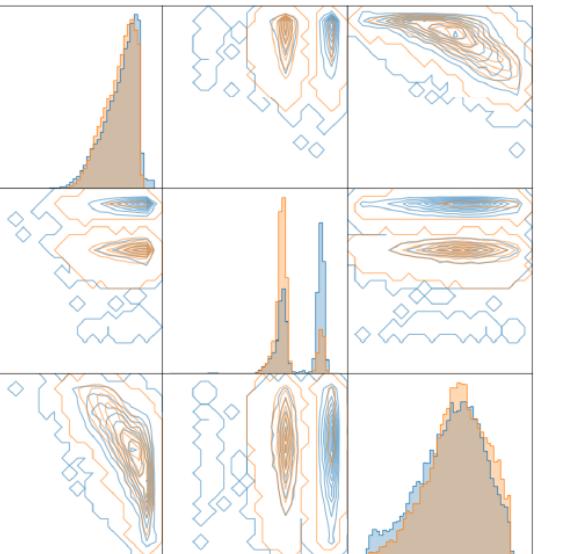
NF



NF+MALA



VAE



VAE+MALA

Method	NF	NF+MALA	VAE	VAE+MALA
KL	0.76	0.59	0.98	0.69

1	2
3	4
5	6
7	8
9	10
11	12
13	14
15	16
17	18
19	20
21	22
23	24
25	26
27	28
29	30
31	32
33	34
35	36
37	38
39	40
41	42
43	44
45	46
47	48
49	50
51	52
53	54
55	56
57	58

Cambridge  
Elements

Non-Local Data Interactions:  
Foundations and Applications

# Generalized Normalizing Flows via Markov Chains

Paul Lyonel Hagemann,  
Johannes Hertrich, and  
Gabriele Steidl

Berlin Mathematics Research Center



Funded under Germany's Excellence Strategy by



Deutsche  
Forschungsgemeinschaft



# Oberwolfach Seminar

## Variational and Information Flows in Machine Learning and Optimal Transport

Organizers: Wuchen Li, Columbia  
Bernhard Schmitzer, Göttingen  
Gabriele Steidl, Berlin  
Francois-Xavier Vialard, Paris  
Date (ID): 19 – 25 November 2023 (2347b)  
Deadline: 1 September 2023

Variational and stochastic flows are now ubiquitous in machine learning and generative modeling. Indeed, many such models can be interpreted as flows from a latent distribution to the sample distribution and training corresponds to finding the right flow vector field. Optimal transport and diffeomorphic flows provide powerful frameworks to analyze such trajectories of distributions with elegant notions from differential geometry, such as geodesics, gradient and Hamiltonian flows. Recently, mean field control and mean field games offer a general optimal control variational problems on the learning problem. How do these tools lead us to a better understanding and further development of machine learning and generative models?

The Oberwolfach Seminar will address the topic from different points of view taking in particular recent developments in machine learning into account. The target audience is PhD students and post-doctoral researchers wishing to be quickly immersed in this modern, active research area. Priority will be given to young, motivated researchers.

Please see the website of the seminar for detailed information:

[www.mfo.de/occasion/2347b](http://www.mfo.de/occasion/2347b)

The seminar takes place at the Mathematisches Forschungsinstitut Oberwolfach. The Institute covers board and lodging. By the support of the Carl Friedrich von Siemens Foundation travel expenses can be reimbursed up to 150 EUR in average per person (against copies of travel receipts). The number of participants is restricted to 25.

**Applications including title, ID and date** of the intended seminar, together with **one pdf-file attached** containing

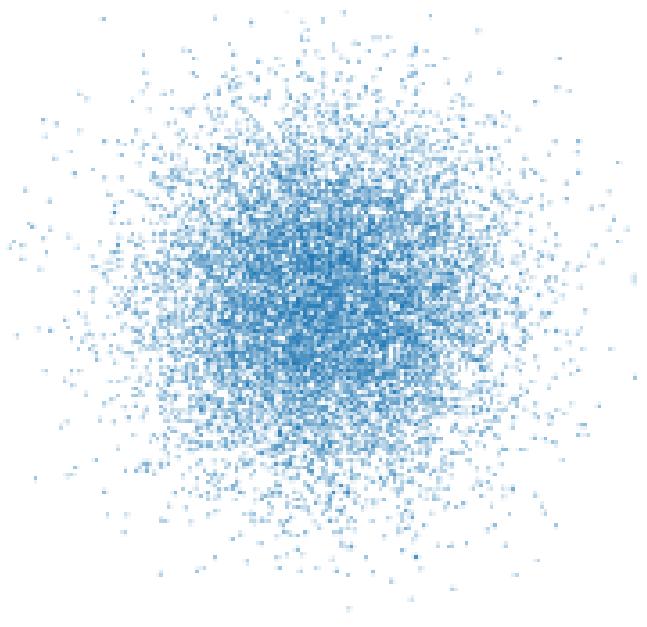
- full name and address, incl. e-mail address
- short CV and publication list
- present position, university
- name of supervisor of Ph.D. thesis
- a short summary of previous work and interest

should be **sent by e-mail** via [seminars@mfo.de](mailto:seminars@mfo.de) until 1 September 2023 to:

Prof. Dr. Matthias Hieber  
Mathematisches Forschungsinstitut Oberwolfach  
Schwarzwaldstr. 9 – 11  
77709 Oberwolfach  
Germany

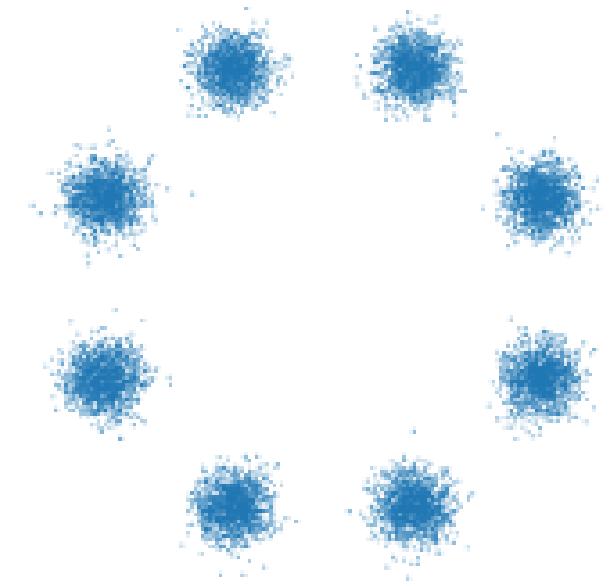




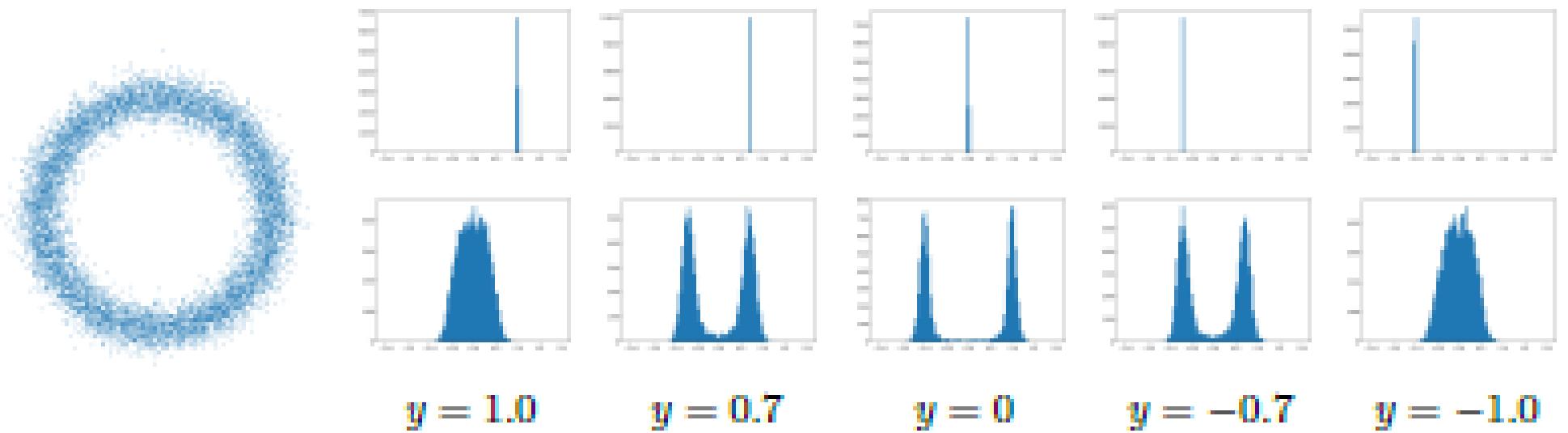


Latent distribution  $P_z$

$$\xrightarrow{\mathcal{T}}$$
  
$$\xleftarrow{\mathcal{T}^{-1}}$$

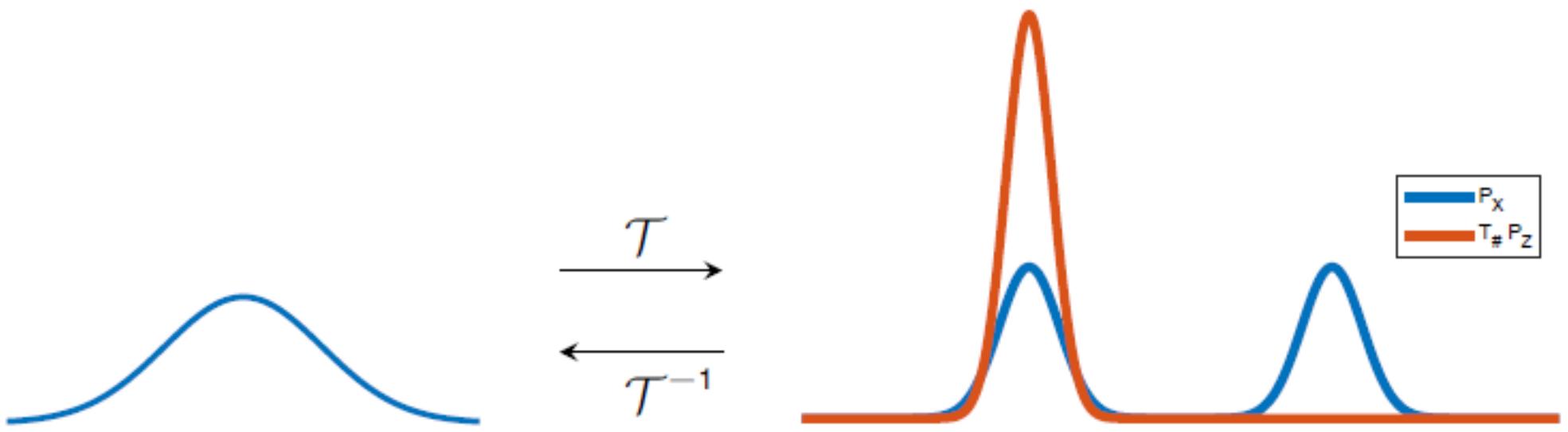


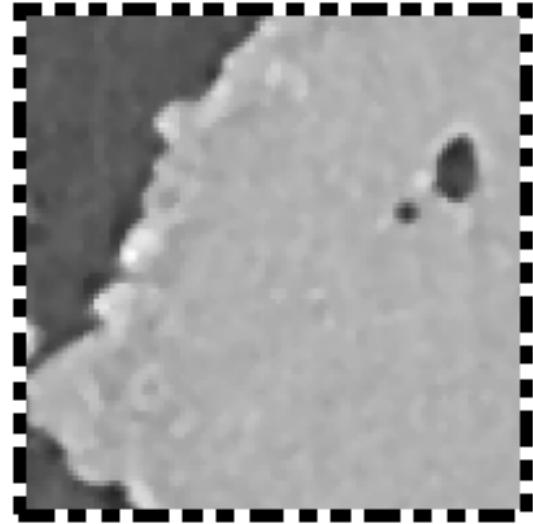
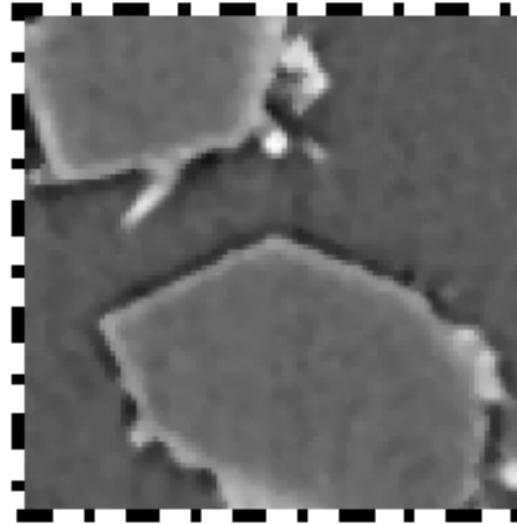
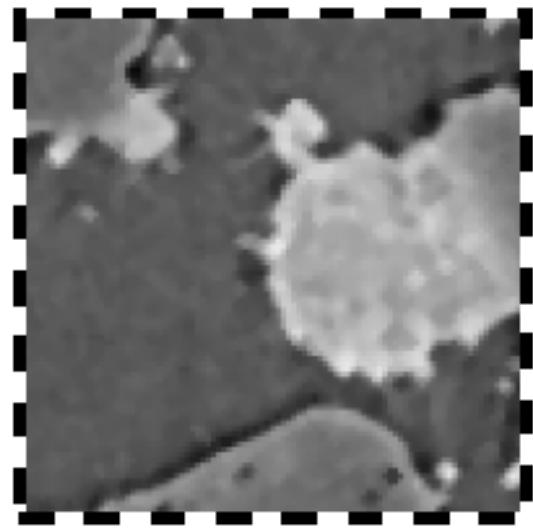
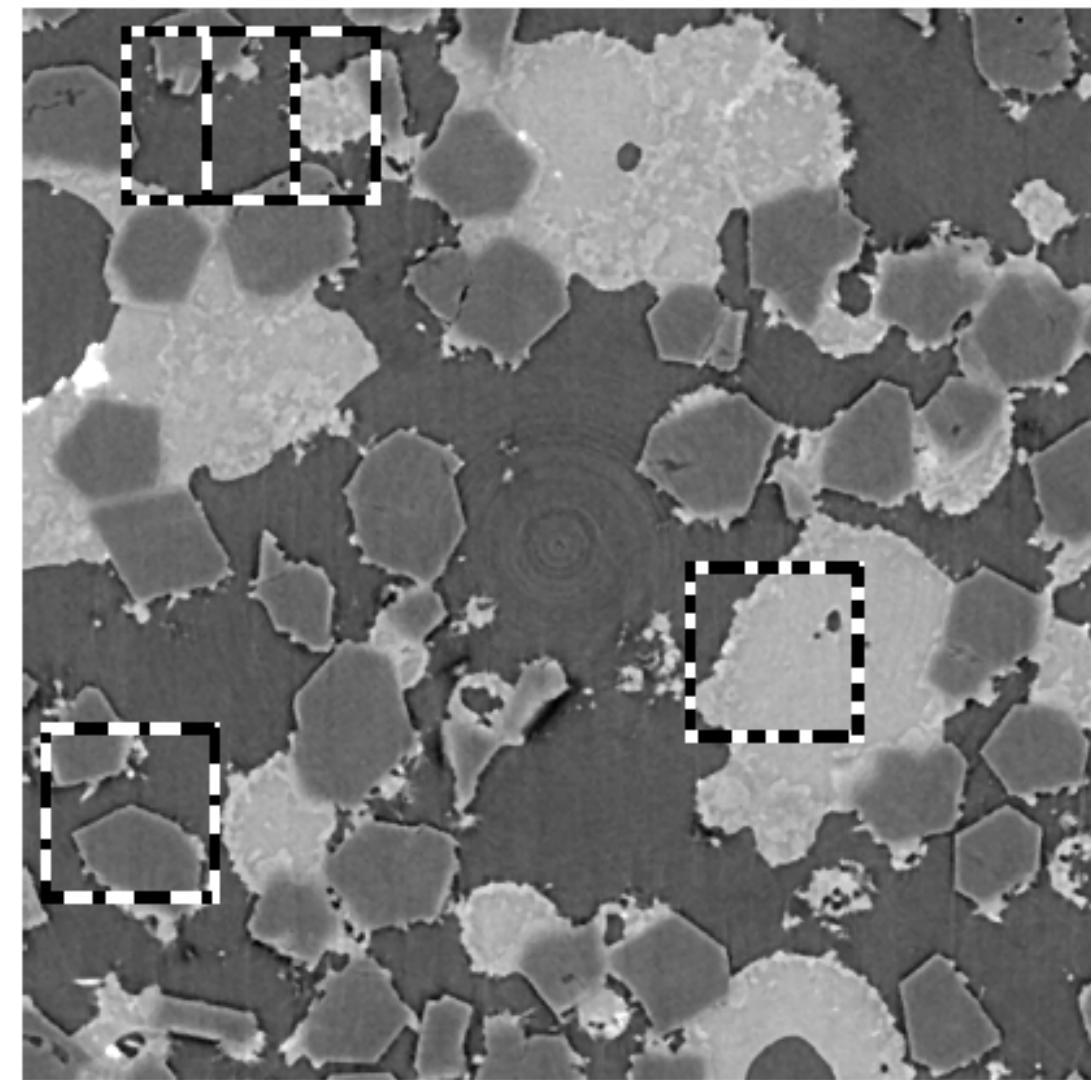
Data distribution  $P_x$

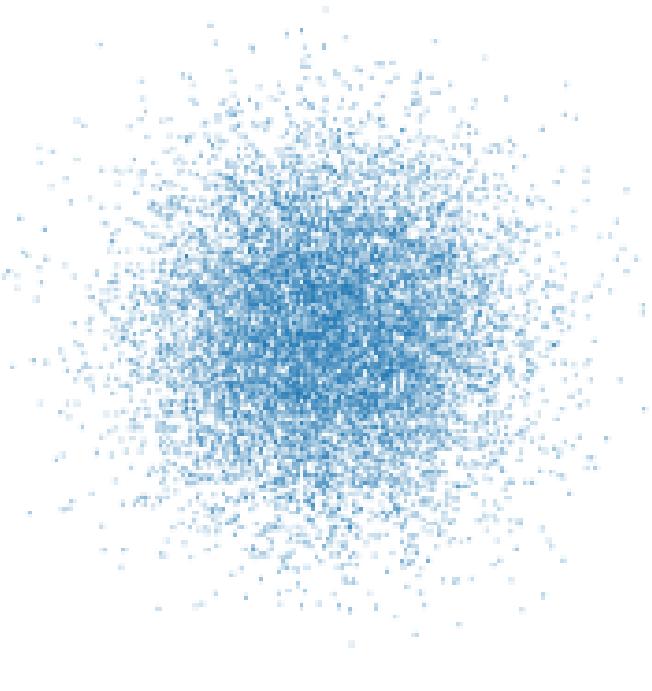


**Figure 2:** Left: Samples from the prior distribution of  $X$  for the circle example. Right: Histograms of samples from the reconstructed posterior distribution  $P_{X|Y=y} \approx \mathcal{T}(\cdot, y)^{-1} P_Z$  for  $y \in \{1, 0.7, 0, -0.7, -1\}$  within the circle example. Top: first coordinate, Bottom: second coordinate.



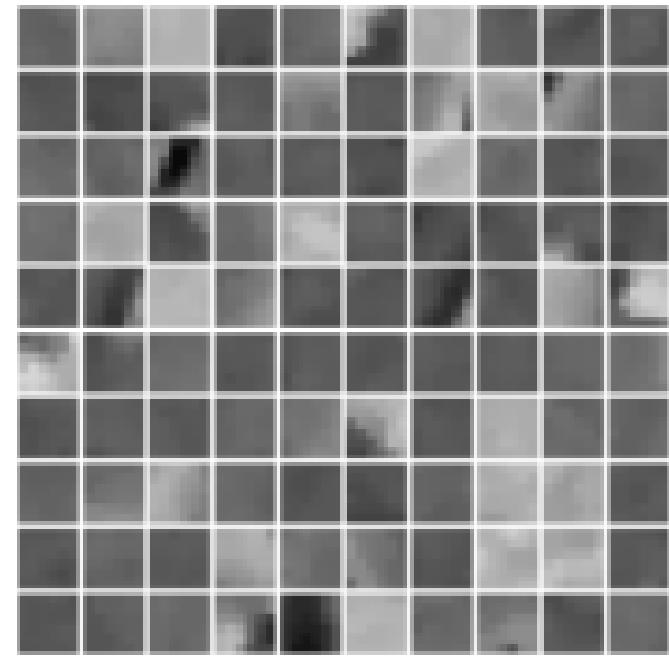




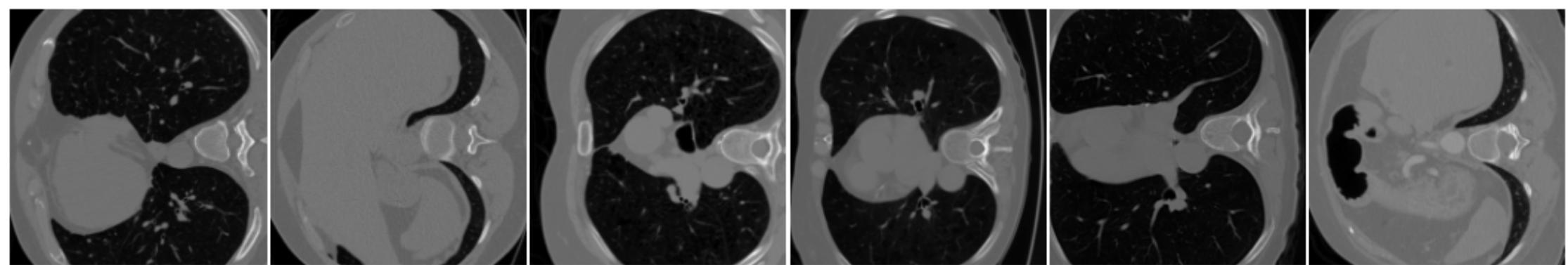


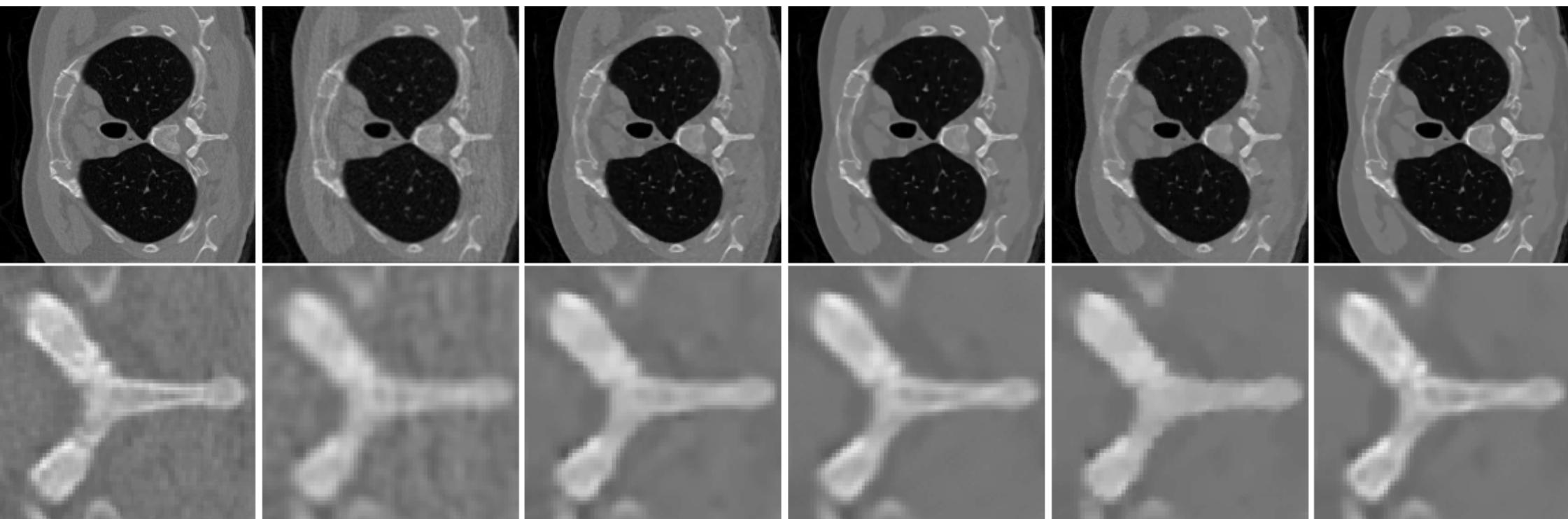
Latent distribution  $P_Z$   
on  $\mathbb{R}^{36}$

$$\begin{array}{c} \xrightarrow{\mathcal{T}} \\ \xleftarrow{\mathcal{T}^{-1}} \end{array}$$



Patch distribution  $P_X$   
on  $\mathbb{R}^{36}$





Ground truth

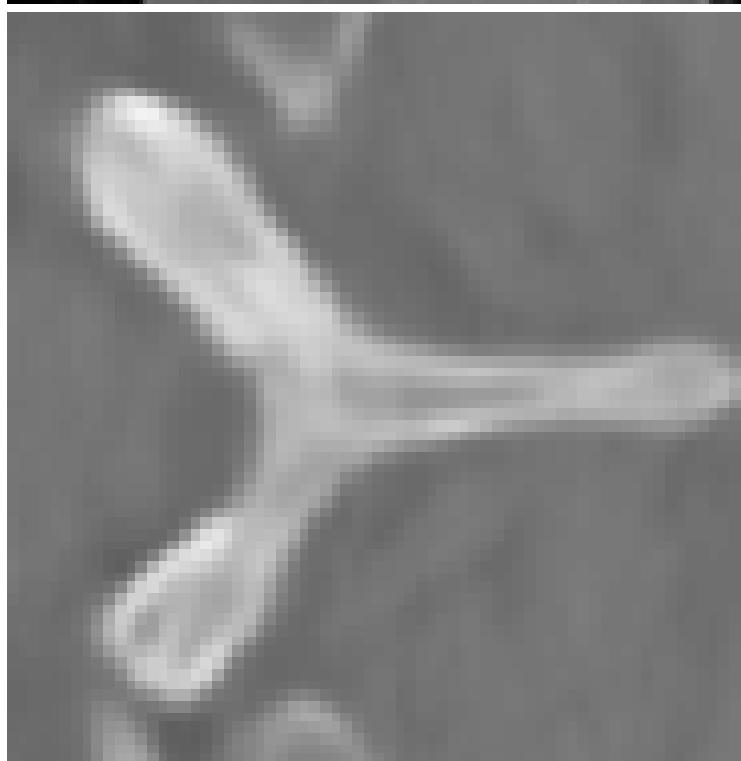
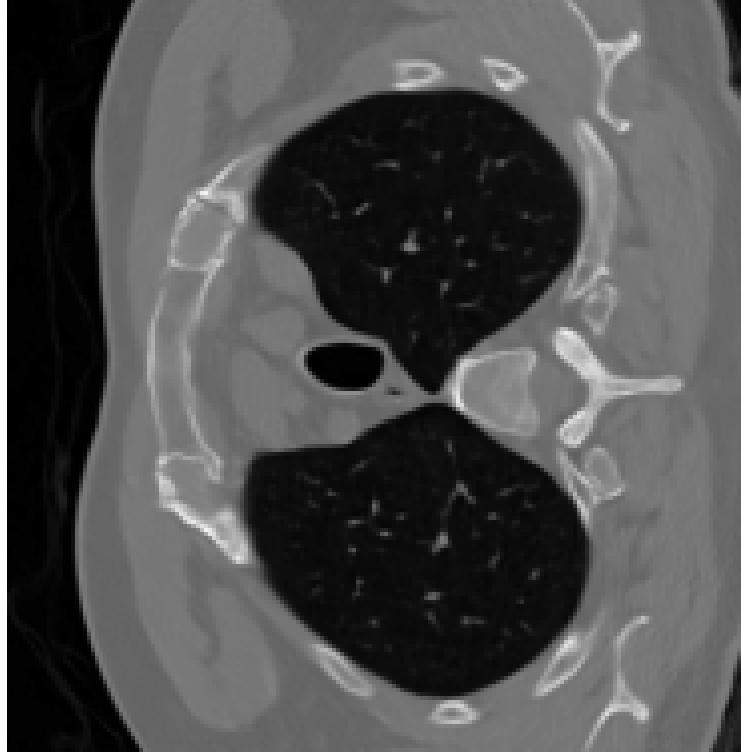
FBP

DIP+TV

EPLL

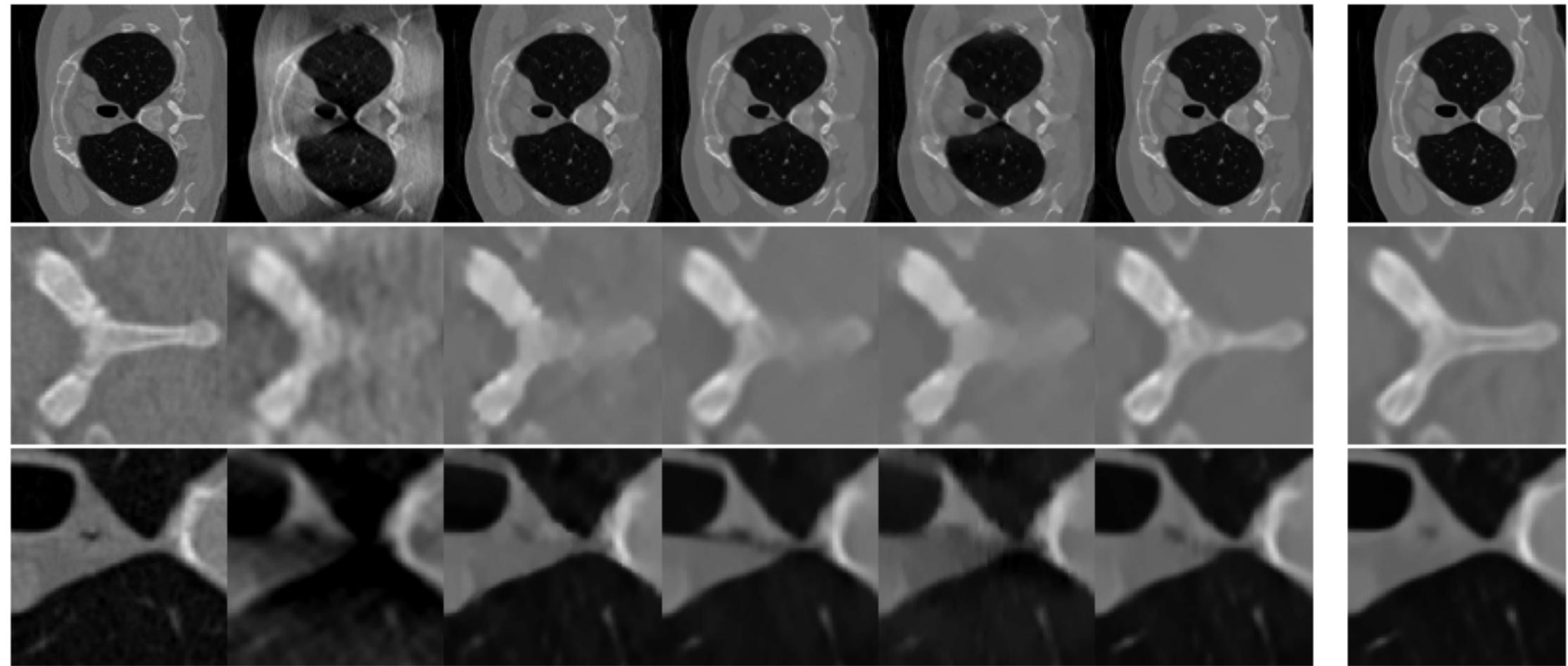
localAR

patchNR



FBP+UNet

	FBP	DIP + TV	EPLL	localAR	patchNR	FBP+UNet (data-based)
PSNR	$30.37 \pm 2.95$	$34.45 \pm 4.20$	$34.89 \pm 4.41$	$33.64 \pm 3.74$	<b><math>35.19 \pm 4.52</math></b>	$35.48 \pm 4.52$
SSIM	$0.739 \pm 0.141$	$0.821 \pm 0.147$	$0.821 \pm 0.154$	$0.807 \pm 0.145$	<b><math>0.829 \pm 0.152</math></b>	$0.837 \pm 0.143$
Runtime	0.03s	1514.33s	36.65s	30.03s	48.39s	0.46s



Ground truth

FBP

DIP+TV

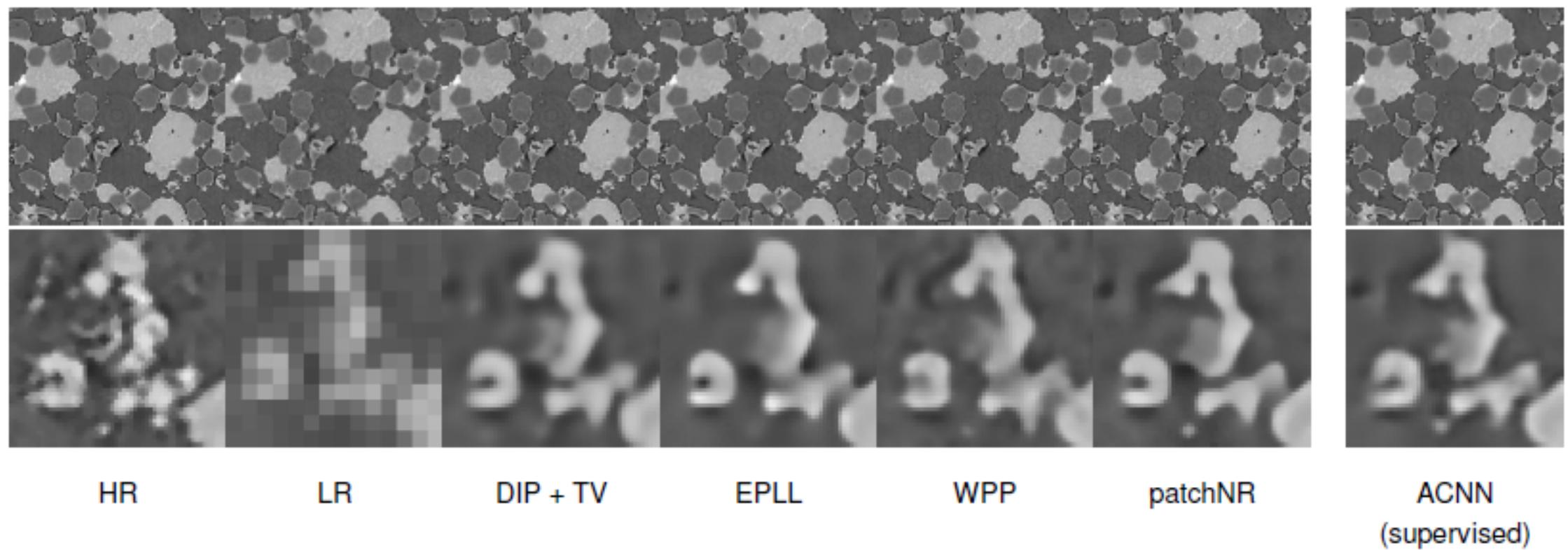
EPLL

localAR

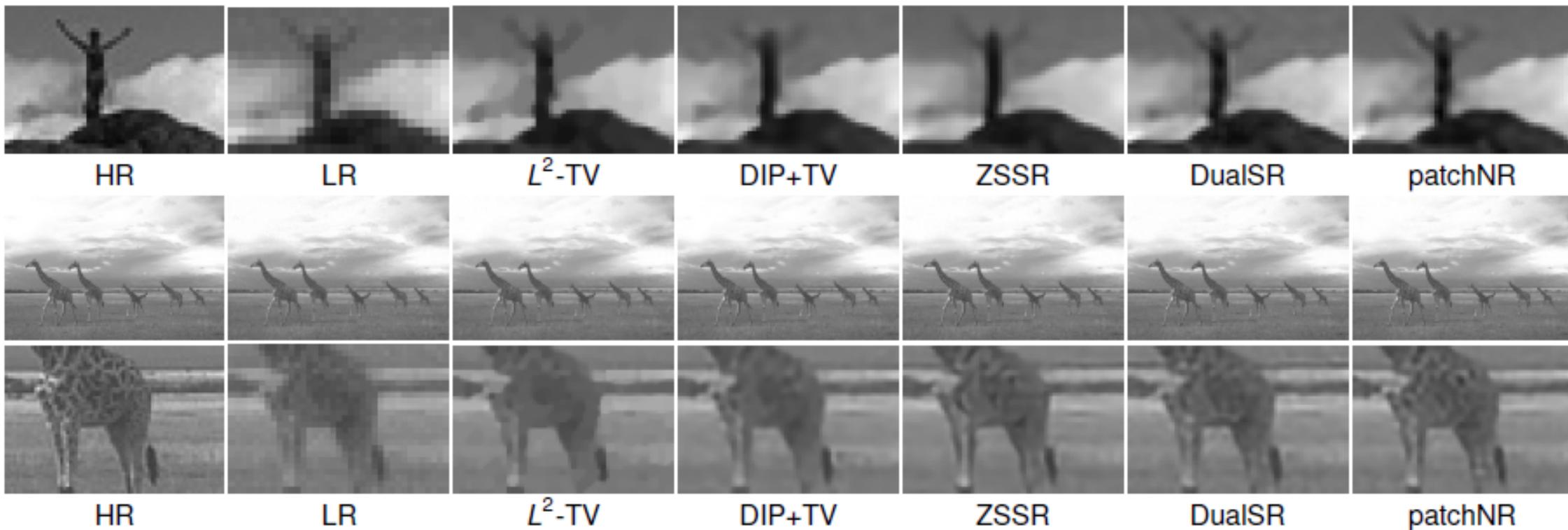
patchNR

FBP+UNet (s)

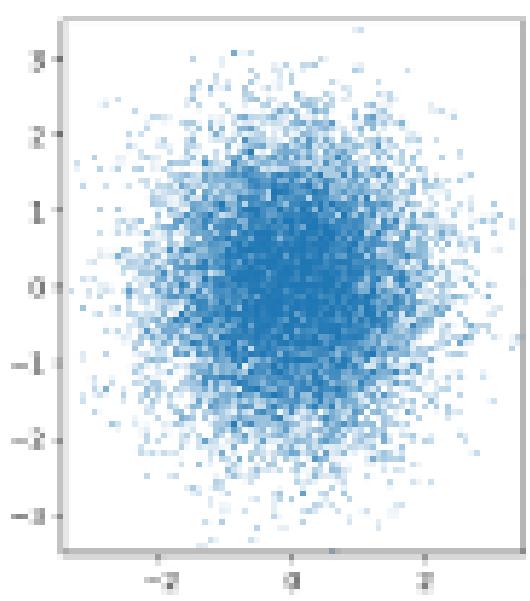
	FBP	DIP + TV	EPLL	localAR	patchNR	FBP+UNet (data-based)
PSNR	$21.96 \pm 2.25$	$32.57 \pm 3.25$	$32.78 \pm 3.46$	$31.06 \pm 2.95$	<b><math>33.20 \pm 3.55</math></b>	$33.75 \pm 3.58$
SSIM	$0.531 \pm 0.097$	$0.803 \pm 0.146$	$0.801 \pm 0.151$	$0.779 \pm 0.142$	<b><math>0.811 \pm 0.151</math></b>	$0.820 \pm 0.140$
Runtime	0.02s	1770.89s	127.21s	53.47s	485.93s	0.53s



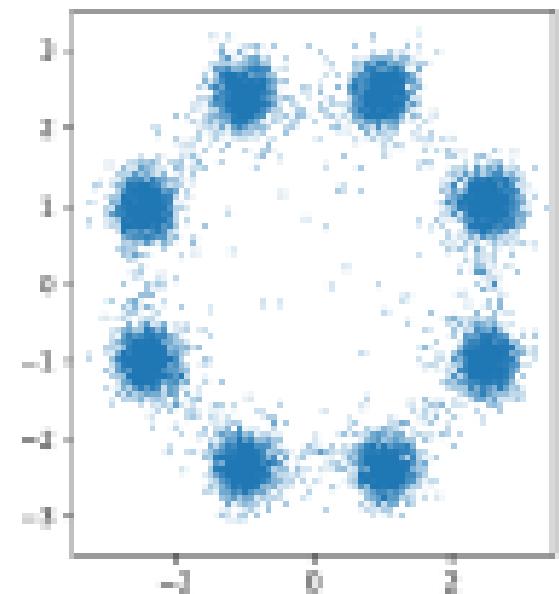
	bicubic (not shown)	DPIR (not shown)	DIP+TV	EPLL	WPP	patchNR	ACNN (data-based)
PSNR	$25.63 \pm 0.56$	$27.78 \pm 0.53$	$27.99 \pm 0.54$	$28.11 \pm 0.55$	$27.80 \pm 0.37$	<b><math>28.53 \pm 0.49</math></b>	$28.89 \pm 0.53$
SSIM	$0.699 \pm 0.012$	$0.770 \pm 0.011$	$0.764 \pm 0.007$	$0.779 \pm 0.010$	$0.749 \pm 0.011$	<b><math>0.780 \pm 0.008</math></b>	$0.804 \pm 0.010$
Runtime	0.0002s	56.62s	234.00s	60.28s	387.28s	150.79s	0.03s



	$L^2$ -TV	DIP+TV	ZSSR	DualSR	patchNR
PSNR	28.35 $\pm$ 3.55	28.44 $\pm$ 3.69	28.83 $\pm$ 3.57	28.64 $\pm$ 3.47	<b>29.08</b> $\pm$ 3.58
SSIM	0.820 $\pm$ 0.072	0.821 $\pm$ 0.087	0.834 $\pm$ 0.066	0.829 $\pm$ 0.061	<b>0.846</b> $\pm$ 0.061
Runtime	13.12s	171.51s	56.64s	53.47s	132.36s



$$\xrightarrow{\quad \mathcal{T} \quad}$$



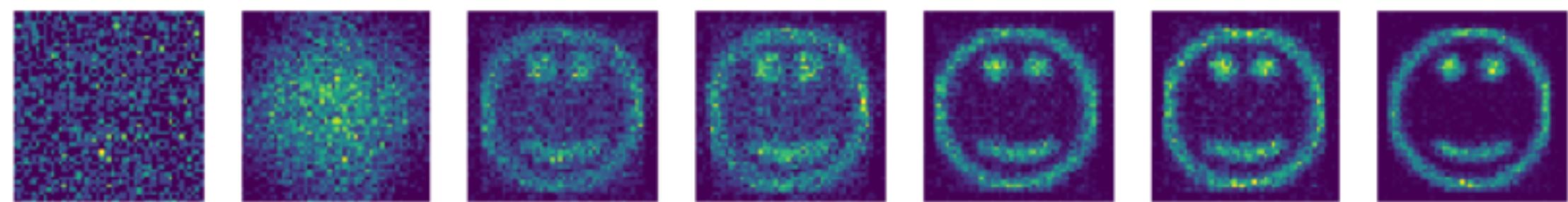
Berlin Mathematics Research Center

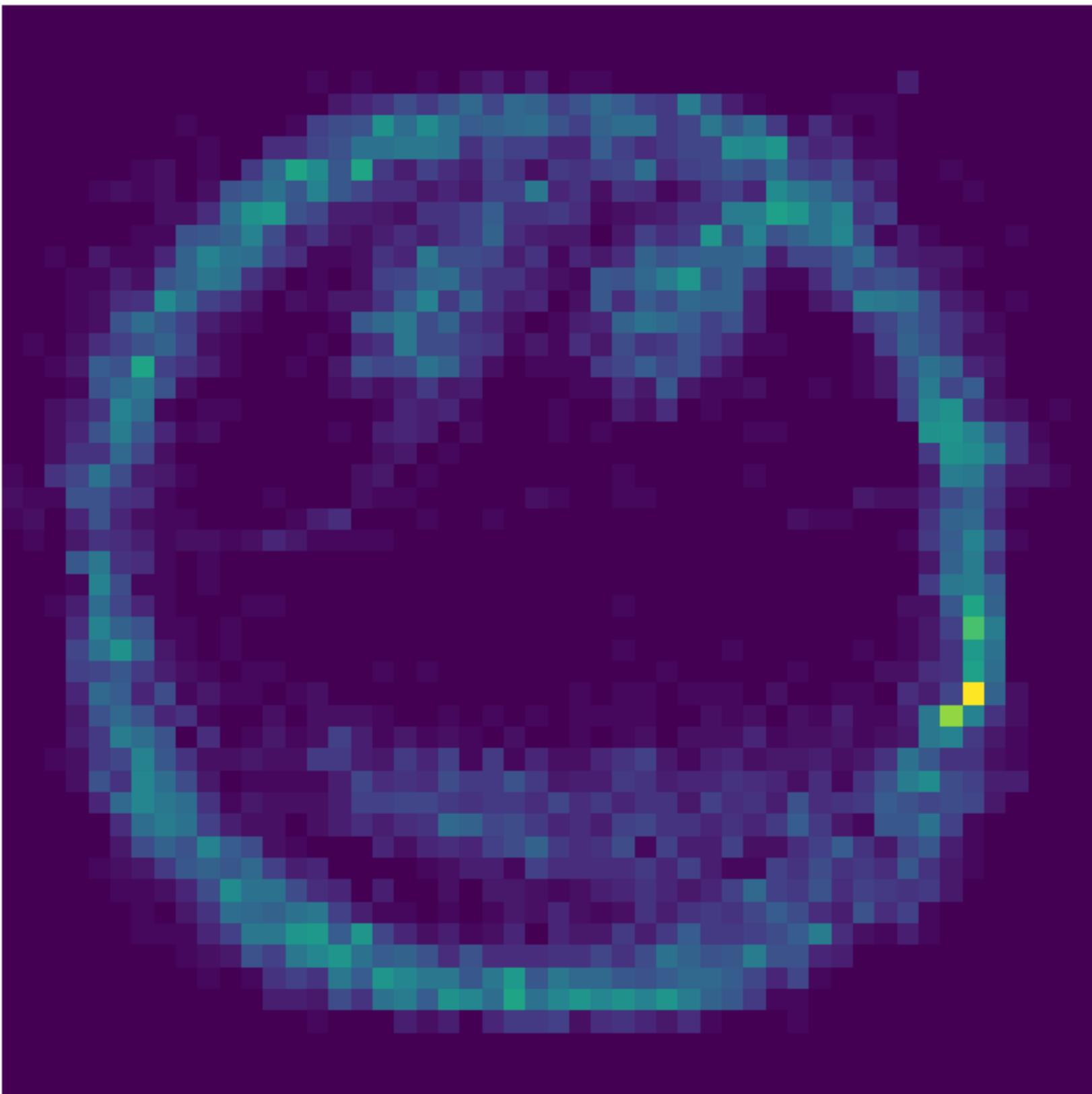


Funded under Germany's Excellence Strategy by



Deutsche  
Forschungsgemeinschaft





### Discrete Measures

$$\begin{array}{c|cccc} & v_1 & \dots & v_n \\ \mu_1 & & & & \\ \vdots & & \pi_{ij} & & \\ \mu_n & & & & \end{array}$$

$$\pi = \sum_{ij} \pi_{ij} \delta_{(x_i, y_j)}$$

$$\pi 1 = \mu$$

$$\pi^T 1 = 0$$

$$M \times K(x_i, y_j) = \pi_{ij}$$

$$M \times K = \pi$$

$$k_{ij} := \frac{\pi_{ij}}{\mu_i} \rightarrow \kappa := (k_{ij})$$

### Markov kernel

$$K(x_i, \cdot) = \sum_e k_{ie} \delta_{y_e}$$

$$P_{Y|x=x_i}$$

$$\rightarrow P_{Y|x=x_i}(y_i) = k_{ii}$$

$$K 1 = 1$$

$$\rightarrow \sum_j \frac{\pi_{ij}}{\mu_i} = \frac{\mu_i}{\mu_i} = 1$$

$$K^T \mu = 0$$

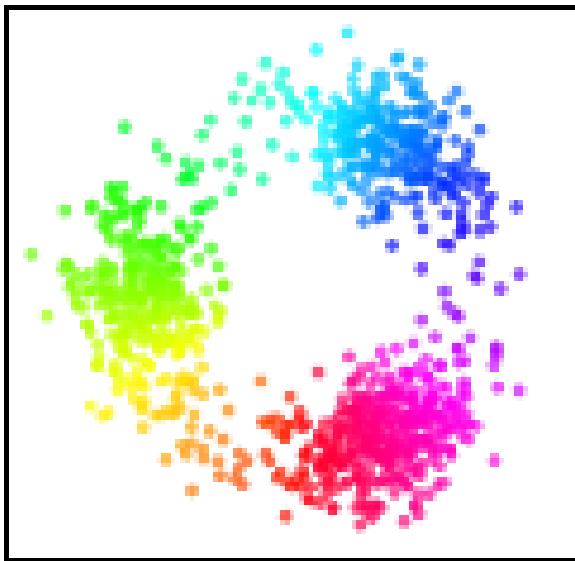
$$\rightarrow \sum_i \frac{\pi_{ij}}{\mu_i} \mu_i = 0_j$$

P stochastic matrix ( $P \geq 0, P^T 1 = 1$ )

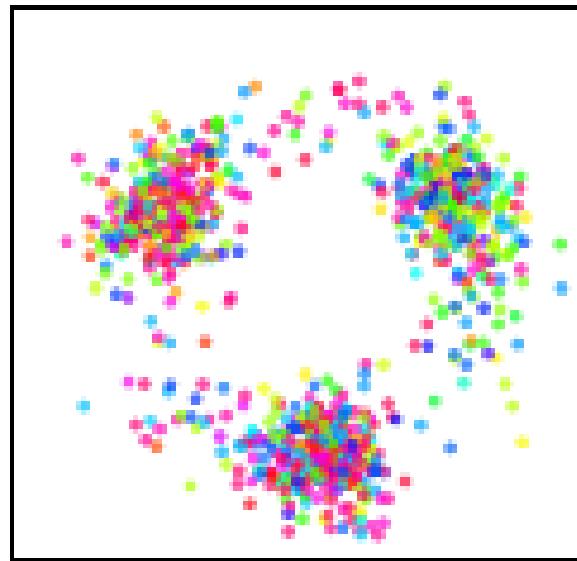
Frobenius-Perron:  $P$  pos. diagonals, mixed.

$\rightarrow P$  has eigenvector  $e$  to largest (simple) eigenvector 1

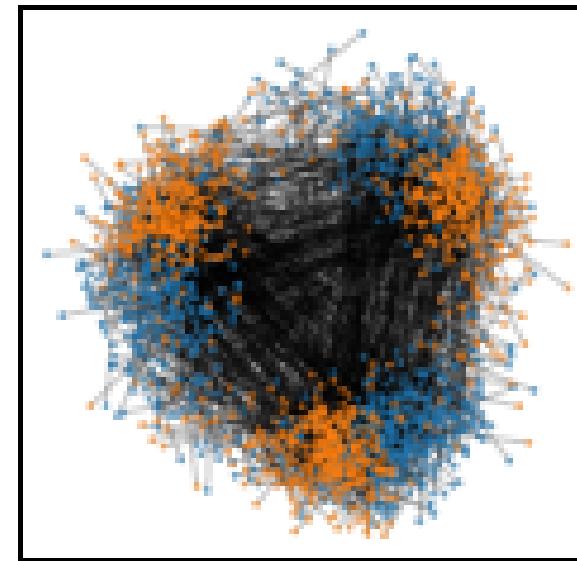
$$\lim_{n \rightarrow \infty} P^n x = e, \|x\|_1 = \|e\|_1$$



(a) Data at time  $t = 0$ .

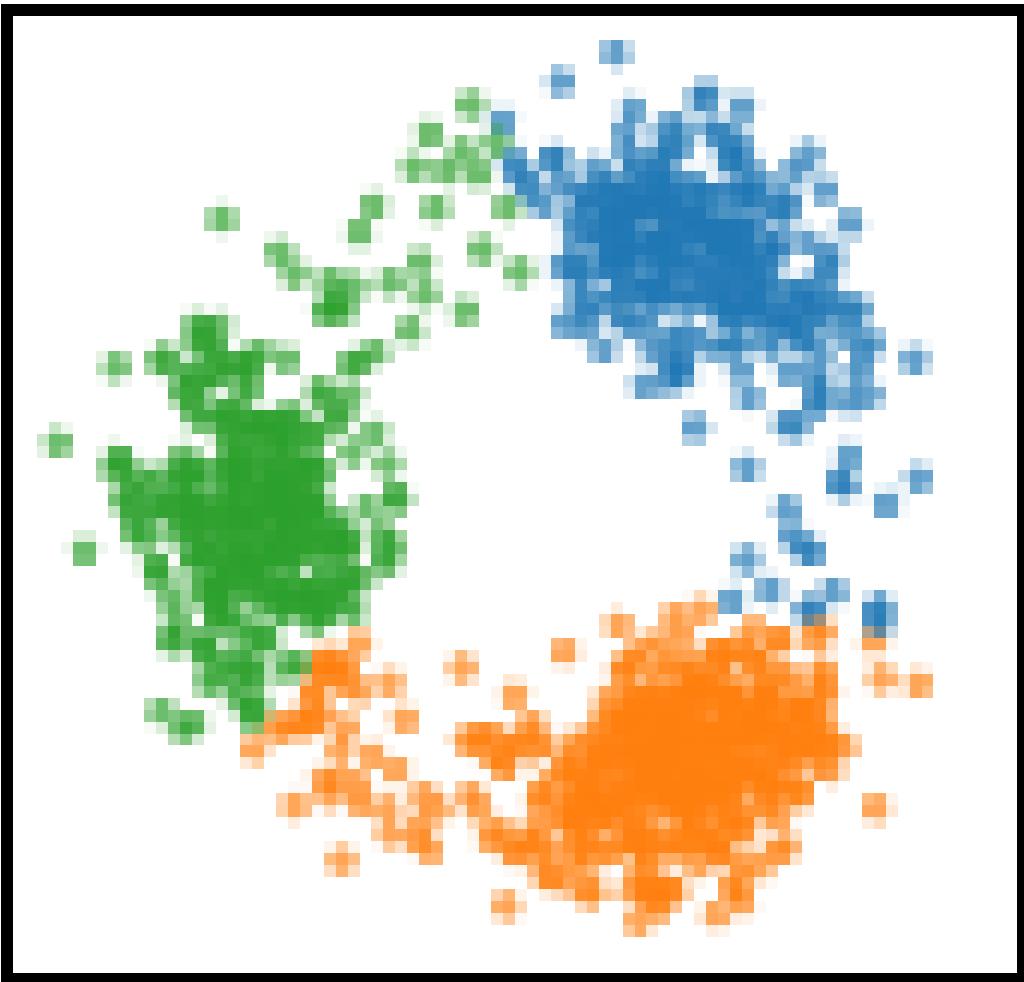


(b) Data at time  $t = 1$ .

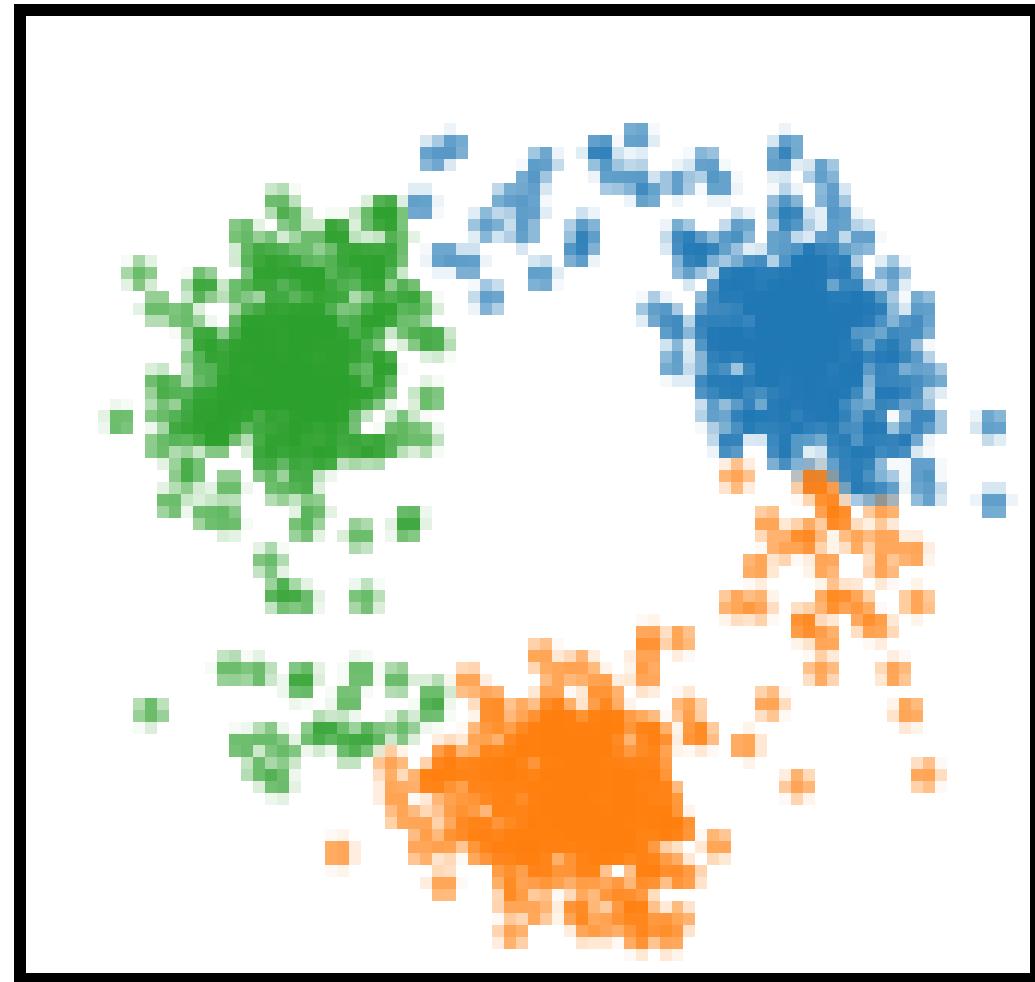


(c) Movement of particles.

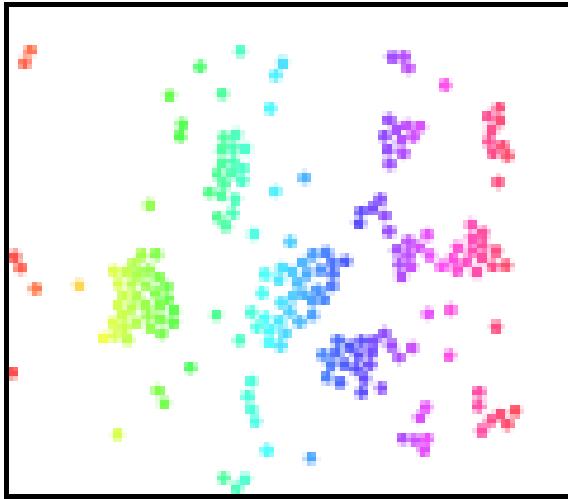
**Figure 7:** Particles moving in a potential with circular driving force. The color scheme illustrates the particle mixing.



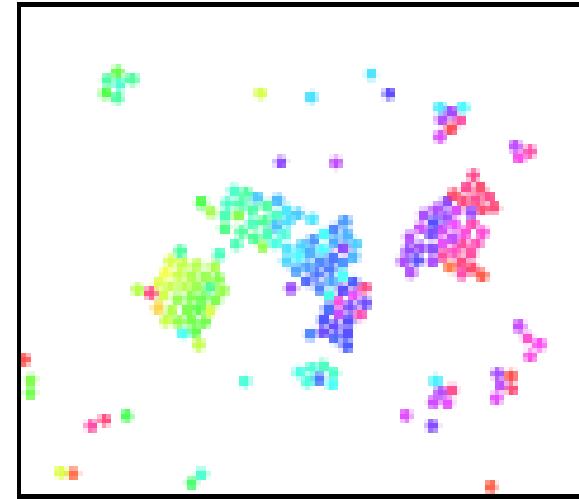
(a) Hand,  $t = 0$ .



(b) Hand,  $t = 1$ .

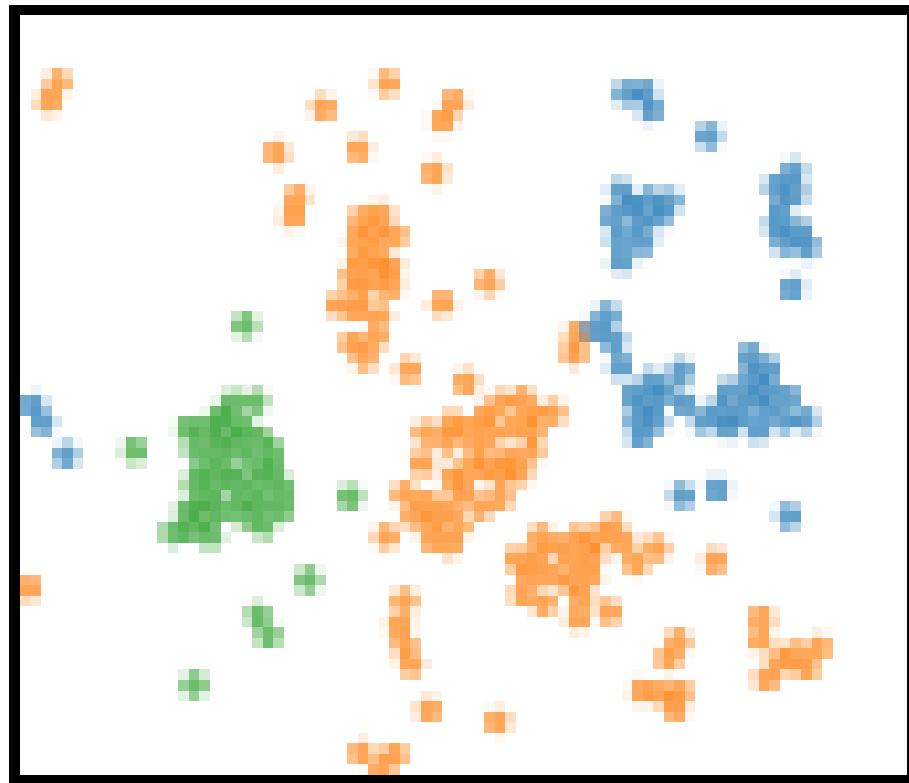


(a) Data at time  $t = 0$ .

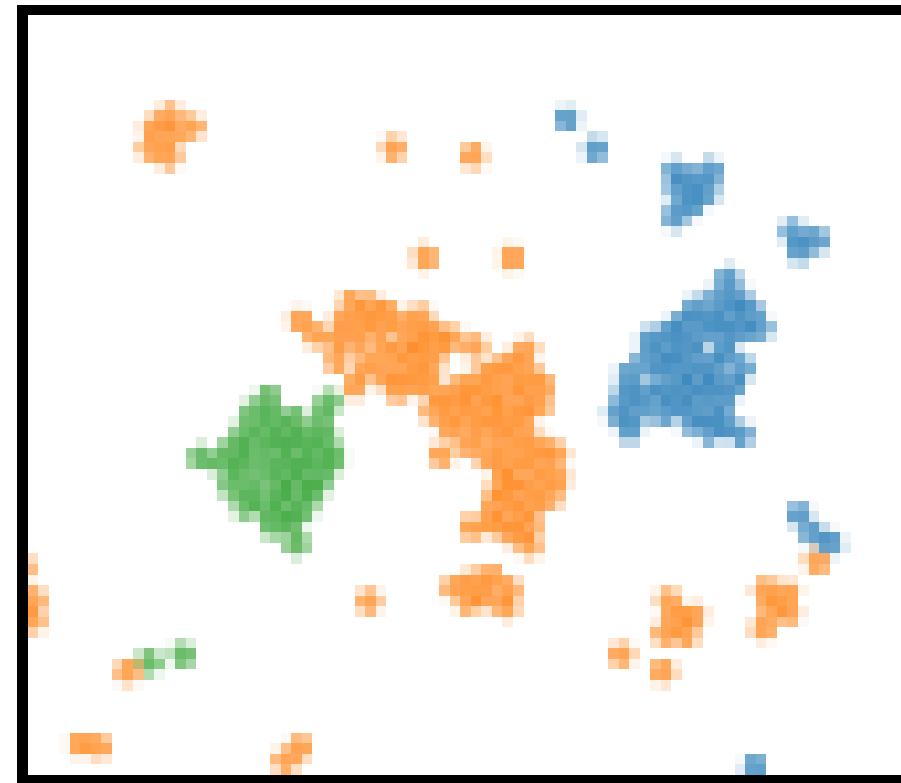


(b) Data at time  $t = 1$ .

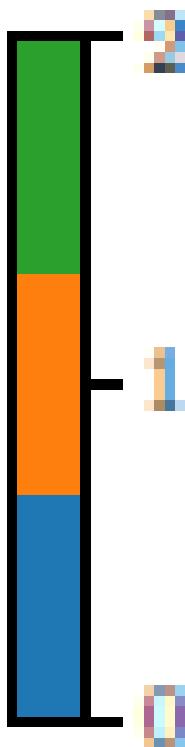
Figure 10: Particles moving in a potential with circular driving force. Again, the color scheme illustrates the particle mixing.



(a) Hand,  $t = 0$ .



(b) Hand,  $t = 1$ .



**Z**

$X_0$

$$X_0 \xrightarrow{\tau_1} X_1 \xrightarrow{\tau_2} \dots$$

$$Y_0 \xleftarrow{\tau_1^{-1}} Y_1 \xleftarrow{\tau_2^{-1}} \dots$$

$$\dots \xrightarrow{\tau_T} X_T$$

$$\dots \xleftarrow{\tau_T^{-1}} Y_T$$

**Y**

**X**

$$\int_{A \times B} p(x, y) dP_{(X, Y)} = \int_A \int_B p(x, y) dP_Y(y) dP_{X|Y=x}(x)$$

$$\textcircled{1} \quad P := \frac{dP(Y_0, \dots, Y_T)}{dP(X_0, \dots, X_T)} = \frac{P_{Y_T}(x_T)}{P_{X_T}(x_T)} \underbrace{\prod_{t=1}^T f_t(x_{t-1}, x_t)}_{\frac{dP_{Y_{t-1}|Y_t=x_t}}{dP_{X_{t-1}|X_t=x_t}}}$$

$$\textcircled{2} \quad KL(\mu, \nu) = \int \log \underbrace{\frac{d\mu}{d\nu}}_{\text{Radon-Nikodym derivative}} d\mu(x) = \mathbb{E}_{X \sim \mu} [\log \frac{d\mu}{d\nu}]$$

$$KL(P_{(Y_0, \dots, Y_T)}, P_{(X_0, \dots, X_T)}) = \mathbb{E}_{(X, Y) \sim P_{(Y_0, \dots, Y_T)}} [\log P_{(X_0, \dots, X_T)}(x, y)]$$

•  $\textcircled{1}$  can be shown by induction using •

Case  $T=1$  ;  $P$  r.h.s in  $\textcircled{1}$

$$\text{To show : } \frac{dP_{(X_0, X_1)}}{dP_{Y_0}(y_1)} = \frac{dP_{(Y_0, Y_1)}}{dP_{X_0}(x_1)}$$

$$\begin{aligned} & \frac{\int_A P_{Y_1}(x_1)}{\int_A P_{X_1}(x_1)} \frac{dP_{Y_0}(y_1)}{\underbrace{dP_{X_0}(x_1)}_{=x_1}} \frac{dP_{(X_0, X_1)}}{(x_0, x_1)} \\ &= \int_A \int_B \frac{dP_{X_1}(x_1)}{dP_{Y_1}(x_1)} \frac{dP_{(X_0 | X_1=x_1)}(x_0)}{dP_{X_0}(x_1)} \\ &= \int_A \int_B dP_{Y_1}(x_1) dP_{Y_0 | Y_1=x_1}(x_0) \\ &= \int_{A \times B} dP_{(Y_0, Y_1)}(x_0, x_1) \end{aligned}$$

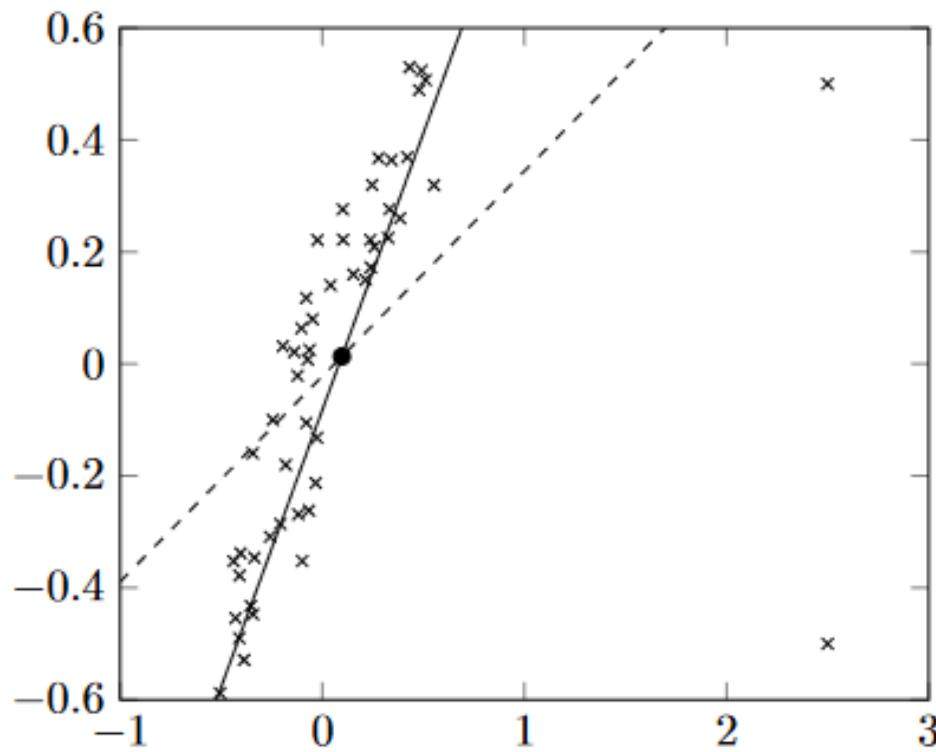
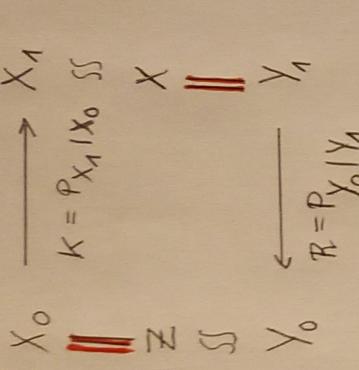


Figure 1: Demonstration of the sensitivity of standard PCA to outliers. The data set consists of 50 points close to a line through the origin and two outliers. The subspace indicated by the dashed line is the result of standard PCA (2), while the solid one corresponds to (4). In both cases the offset  $b$  was chosen as the mean (black dot).

Markov chain ( $x$  layer) :



Aim to learn  $\theta$  such that

$$P_X(x) \approx \int_{\mathbb{R}^n} p_\theta(x|z) p(z) dz$$

$$P_X(A) \approx \int_{\mathbb{R}^n} \kappa(z, A) dP_Z(z)$$

Sample for  $P_X$

- ① Sample from  $Z$
- ② Sample from  $\kappa(z, \cdot)$

evidence

Loss function: Idea:  $E_{X \sim P_X} \left[ \log \underbrace{\frac{p_\theta(x)}{\int_{\mathbb{R}^n} p_\theta(x|z) p(z) dz}}_{=} \right] \rightarrow \max$

How to compare

Approximation of evidence from below: evidence lower bound (ELBO)

$$\begin{aligned} \log(p_\theta(x)) &= E_{Z^n q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x) p_\theta(z|x)}{p_\theta(z|x)} \right] \\ &= E_{Z^n q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x) p_\theta(z|x)}{q_\phi(z|x)} \right] + E_{Z^n q_\phi(\cdot|x)} \left[ \log \frac{q_\phi(z|x)}{p_\theta(z|x)} \right] \\ &\stackrel{\text{Bayes}}{=} E_{Z^n q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x|z) p_\theta(z)}{q_\phi(z|x)} \right] + KL \left( q_\phi(\cdot|x), p_\theta(\cdot|x) \right) \\ &\geq E_{Z^n q_\phi(\cdot|x)} \left[ \log \frac{p_\theta(x|z) p_\theta(z)}{q_\phi(z|x)} \right] \end{aligned}$$

Relation to Markov chains:

$$L_{SNF}(\theta, \varphi) = KL(P(y_0, y_1) \| P(x_0, x_1))$$

$$= E_{\substack{(x_0, x_1) \sim P(y_0, y_1) \\ z \sim X}} \left[ \log \frac{p_{y_1}(x_1)}{p_{x_1}(x_1)} f_1(x_0, x_1) \right]$$

$$L_{SNF}(\theta, \varphi) = E_{(z, x) \sim P(y_0, x)} \left[ \log \frac{p_x(x)}{p_{x_1}(x)} f_1(z, x) \right]$$

$$\begin{aligned}
 f_1(z|x) &= \frac{d P_{Y_0|X=x}}{d P_{Z|X_1=x}}(z) \stackrel{\text{Bayes}}{=} \frac{q_\phi(z|x)}{P_{Z|X_1=x}} = \\
 &= \frac{q_\phi(z|x)}{p_\theta(x|z)} \frac{p_{X_1}(x)}{p_Z(z)}
 \end{aligned}$$

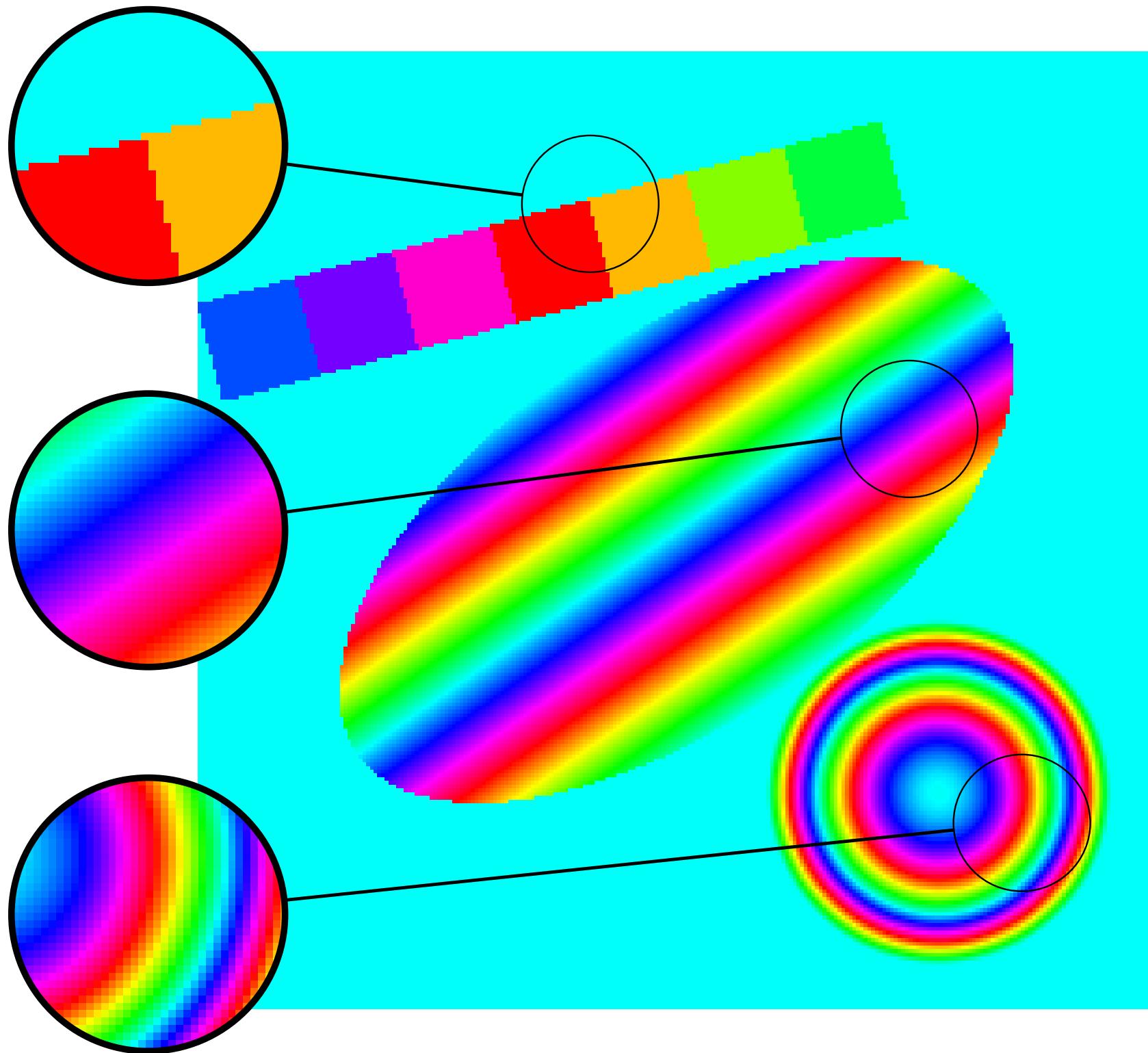
Thus

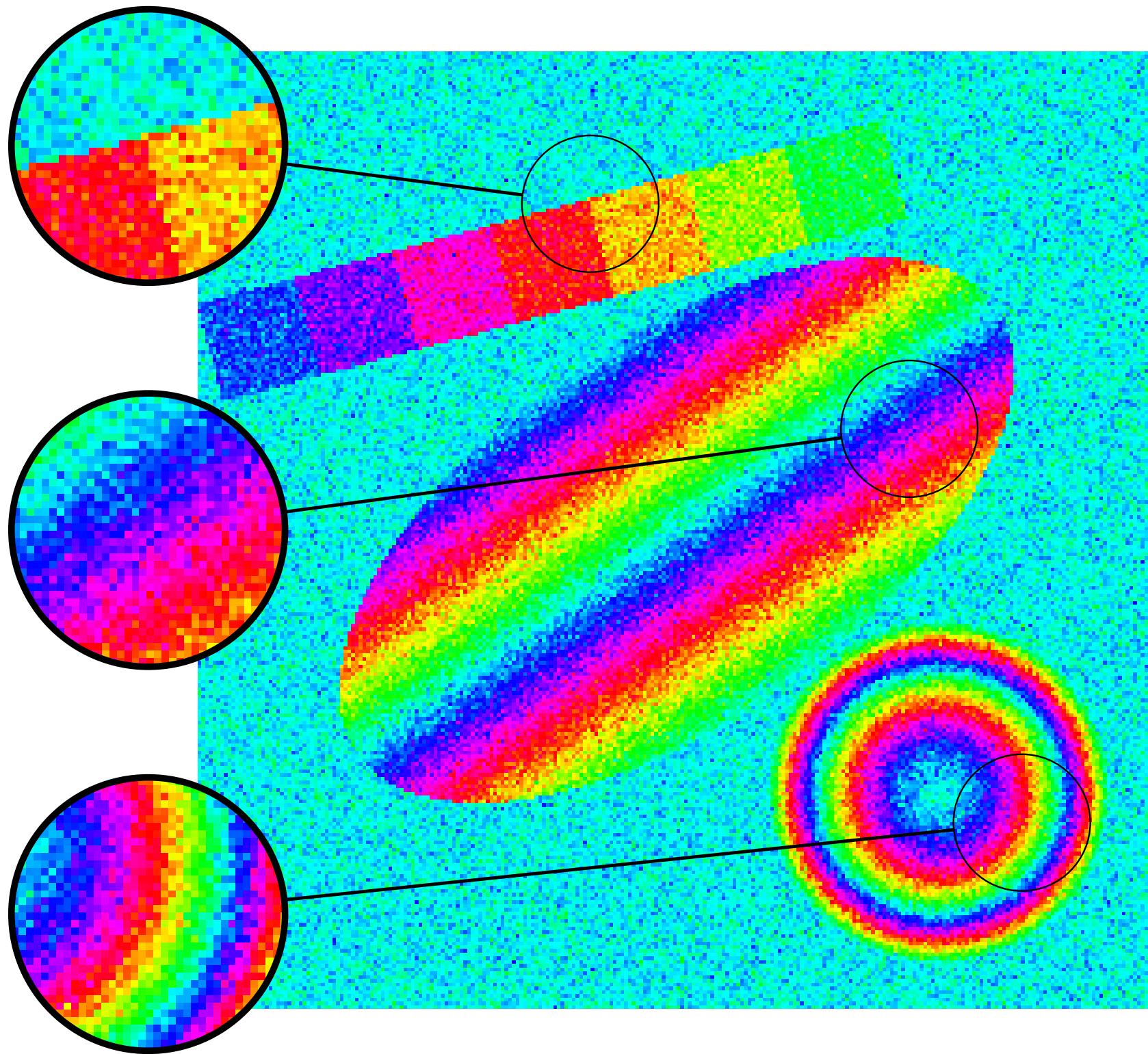
$$L_{SNF}(\theta, \varphi) = E_{(Z, X) \sim P(Y_0, X)} \left[ \log \frac{q_\phi(z|x)}{p_\theta(x|z)} \right]$$

$$\begin{aligned}
 P_{(Y_0, X)} &= P_X \times \underbrace{P_{Y_0|X}_{\mathcal{R}(X, \cdot)}}_{\text{Markov kernel property}} \\
 &= \mathbb{E}_{X \sim P_X} \left( \mathbb{E}_{Z \sim q_\phi(\cdot, x)} \left[ \log \frac{q_\phi(z|x)}{p_\theta(x|z)} \frac{p_X(x)}{p_Z(z)} \right] \right) \\
 &= \mathbb{E}_{X \sim P_X} \left( \mathbb{E}_{Z \sim q_\phi(\cdot, x)} \left[ \log \frac{q_\phi(z|x)}{p_\theta(x|z)} \frac{p_X(x)}{p_Z(z)} \right] \right) \\
 &\quad - \mathbb{E}_{X \sim P_X} \left[ p_X(x) \right] \\
 &\quad \text{Const}
 \end{aligned}$$

$$\begin{aligned}
 \underset{\theta, \varphi}{\operatorname{argmin}} L_{SNF}(\theta, \varphi) &= \underset{\theta, \varphi}{\operatorname{argmax}} E_{X \sim P_X} \left( \mathbb{E}_{Z \sim q_\phi(\cdot, x)} \right)
 \end{aligned}$$

$\Rightarrow$





$-\pi$  $0$  $\pi$ 