

Adaptive construction of Transport Maps for efficient sampling

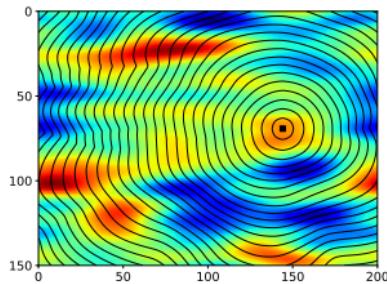
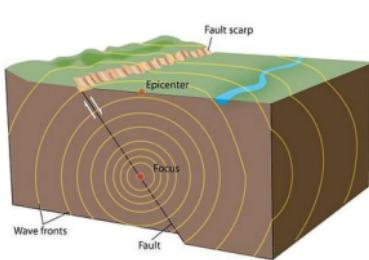
D. Bigoni (**dabi@mit.edu**), A. Spantini, Y.M. Marzouk
Massachusetts Institute of Technology

Past and present contributors:

Tarek El Moselhy, Matthew Parno, Xun Huan, Rebecca Morrison,
Ricardo M. Batista, Benjamin Zhang, Zheng Wang

MCQMC 2018
Rennes – July 5, 2018

Bayesian inference – an oversimplified example



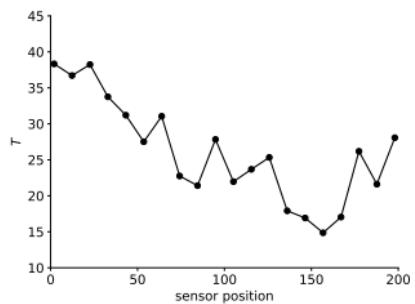
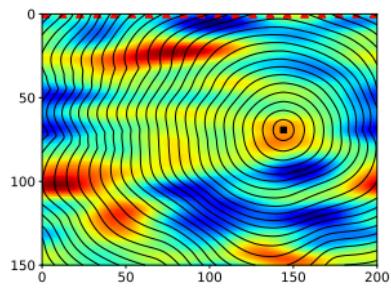
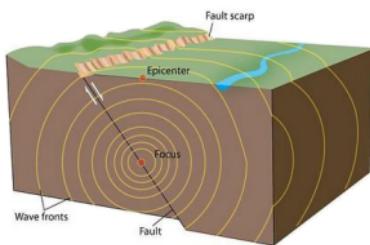
Mathematical model

travel time

$$|\nabla \tilde{G}(\mathbf{x})| = v(\mathbf{x})^{-1}$$

velocity field

Bayesian inference – an oversimplified example



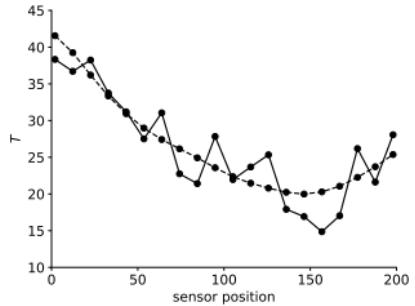
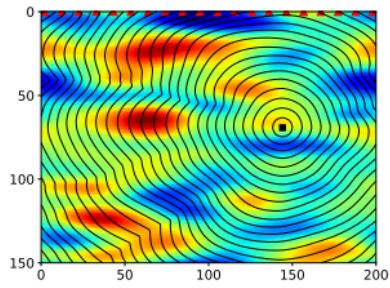
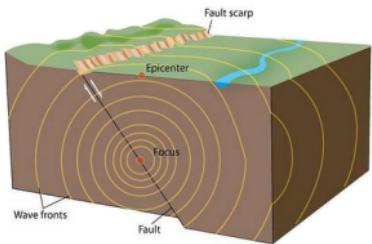
Mathematical model

$$\underbrace{|\nabla G(\mathbf{x})|}_{\text{travel time}} = \underbrace{v(\mathbf{x})^{-1}}_{\text{velocity field}}$$

Observational model

$$\underbrace{\mathbf{d}}_{\text{data}} = \mathbf{G}(\mathbf{v}) + \underbrace{\boldsymbol{\varepsilon}}_{\text{noise}}$$

Bayesian inference – an oversimplified example



Mathematical model

$$\text{travel time} \quad |\nabla G(\mathbf{x})| = v(\mathbf{x})^{-1}$$

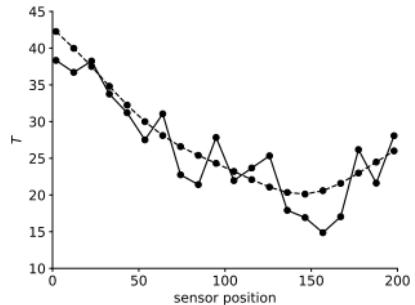
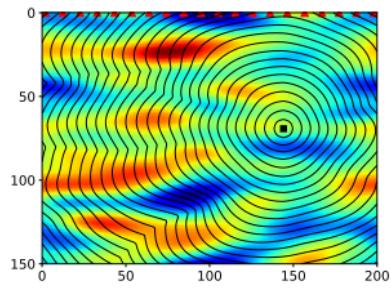
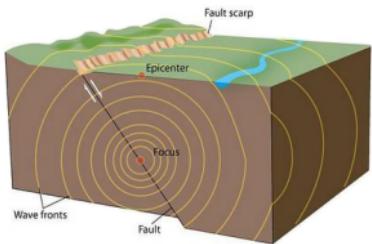
velocity field

Observational model

$$\text{data} \quad \mathbf{d} = \mathbf{G}(\mathbf{v}) + \varepsilon$$

noise

Bayesian inference – an oversimplified example



Mathematical model

$$\text{travel time} \quad |\nabla G(\mathbf{x})| = v(\mathbf{x})^{-1}$$

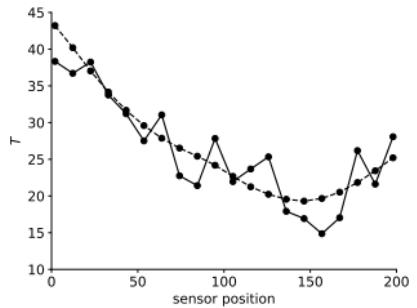
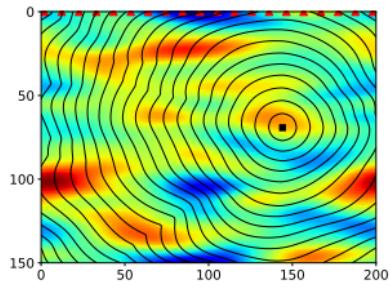
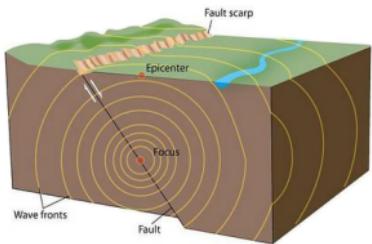
velocity field

Observational model

$$\text{data} \quad \mathbf{d} = \mathbf{G}(\mathbf{v}) + \varepsilon$$

noise

Bayesian inference – an oversimplified example



Mathematical model

$$\text{travel time } |\nabla G(\mathbf{x})| = v(\mathbf{x})^{-1}$$

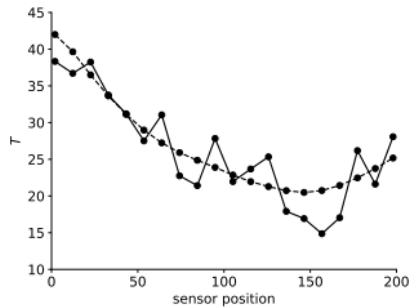
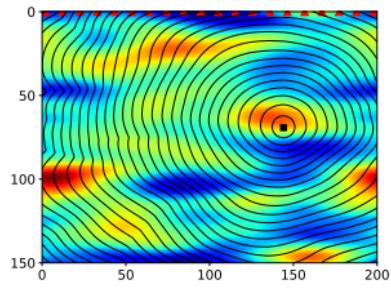
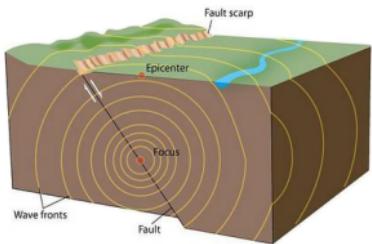
velocity field

Observational model

$$\text{data } \mathbf{d} = \mathbf{G}(\mathbf{v}) + \varepsilon$$

noise

Bayesian inference – an oversimplified example



Mathematical model

$$\text{travel time} \quad |\nabla G(\mathbf{x})| = v(\mathbf{x})^{-1}$$

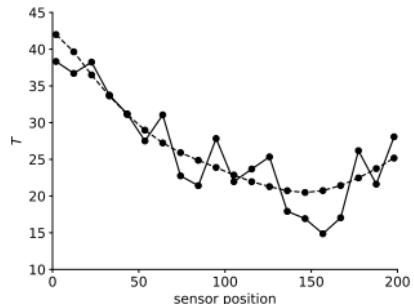
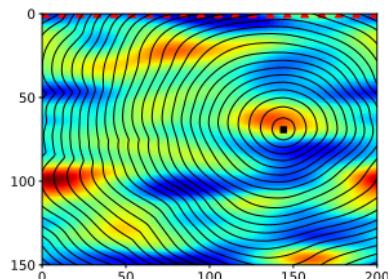
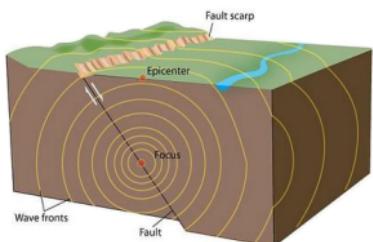
velocity field

Observational model

$$\text{data} \quad \mathbf{d} = \mathbf{G}(\mathbf{v}) + \varepsilon$$

noise

Bayesian inference – an oversimplified example



Mathematical model

$$|\nabla G(\mathbf{x})| = v(\mathbf{x})^{-1}$$

travel time
velocity field

Observational model

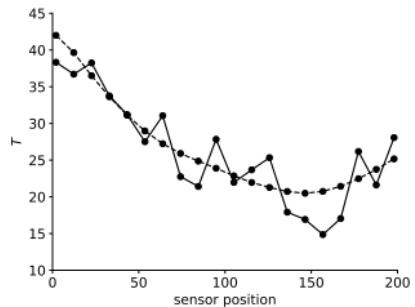
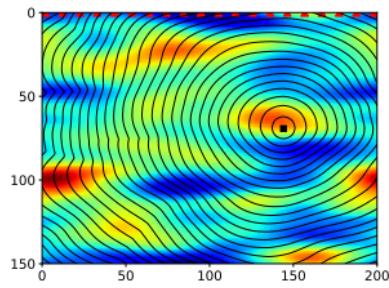
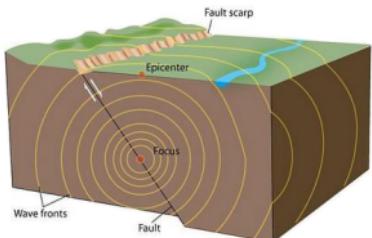
$$\mathbf{d} = \mathbf{G}(\mathbf{v}) + \boldsymbol{\varepsilon}$$

data
noise

Bayesian inference model

$$\underbrace{\pi_{\text{pos}}(\mathbf{v} | \mathbf{d})}_{\text{posterior}} \propto \underbrace{\mathcal{L}_{\mathbf{d}}(\mathbf{v})}_{\text{likelihood}} \underbrace{\pi_{\text{pr}}(\mathbf{v})}_{\text{prior}} = \pi_{\boldsymbol{\varepsilon}}(\mathbf{d} - \mathbf{G}(\mathbf{v})) \pi_{\text{pr}}(\mathbf{v})$$

Bayesian inference – an oversimplified example



Bayesian inference model

$$\underbrace{\pi_{\text{pos}}(\mathbf{v}|\mathbf{d})}_{\text{posterior}} \propto \overbrace{\mathcal{L}_{\mathbf{d}}(\mathbf{v})}^{\text{likelihood}} \underbrace{\pi_{\text{pr}}(\mathbf{v})}_{\text{prior}} = \pi_{\boldsymbol{\varepsilon}}(\mathbf{d} - \mathbf{G}(\mathbf{v})) \pi_{\text{pr}}(\mathbf{v})$$

Decisions under uncertainty

$$\min_{\delta} \int L(\mathbf{v}, \delta) \pi_{\text{pos}}(\mathbf{v}|\mathbf{d}) d\mathbf{v}$$

Goal: characterize $\pi_{\text{pos}}(\mathbf{v}|\mathbf{d})$, i.e.

- construct approximations

$$\int f(\mathbf{v}) \pi_{\text{pos}}(\mathbf{v}|\mathbf{d}) d\mathbf{v} \approx \int f(\mathbf{v}) \tilde{\pi}_{\text{pos}}(\mathbf{v}|\mathbf{d}) d\mathbf{v} \approx \sum_{i=1}^n f(\mathbf{v}^{(i)}) \mathbf{w}^{(i)}$$

- control the error between $\pi_{\text{pos}}(\mathbf{v}|\mathbf{d})$ and $\tilde{\pi}_{\text{pos}}(\mathbf{v}|\mathbf{d})$

Difficulties:

- $\mathbf{v} \in \mathbb{R}^d$ where $d \gg 1$
- The model $\mathbf{G}(\mathbf{v})$ is non-linear
- Evaluation of the model $\mathbf{G}(\mathbf{v})$ is expensive

Outline

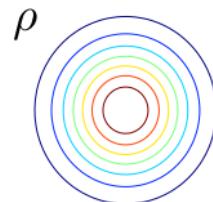
Transport maps

Adaptivity

Off-line learning / On-line inference

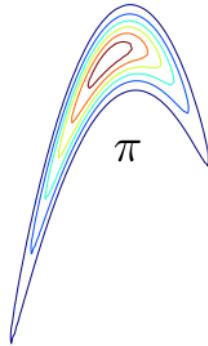
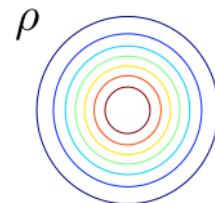
Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$



Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$

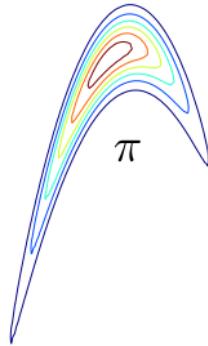
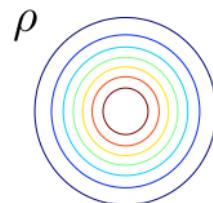


Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

$$\text{PF} \quad T_\sharp \rho = \rho \circ T^{-1} |\nabla T^{-1}|$$

$$\text{PB} \quad T^\sharp \pi = \pi \circ T |\nabla T|$$



Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

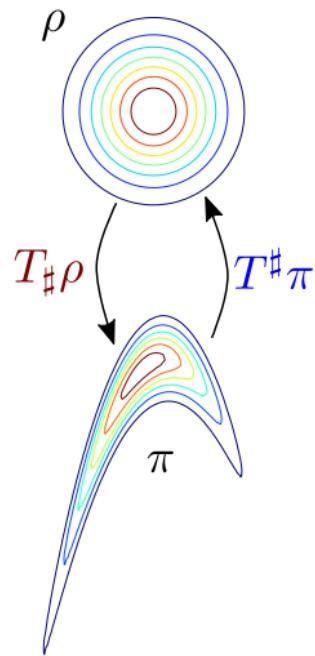
$$\text{PF} \quad T_\sharp \rho = \rho \circ T^{-1} |\nabla T^{-1}|$$

$$\text{PB} \quad T^\sharp \pi = \pi \circ T |\nabla T|$$

- We want T such that

$$\text{PF} \quad T_\sharp \rho = \pi$$

$$\text{PB} \quad T^\sharp \pi = \rho$$



Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

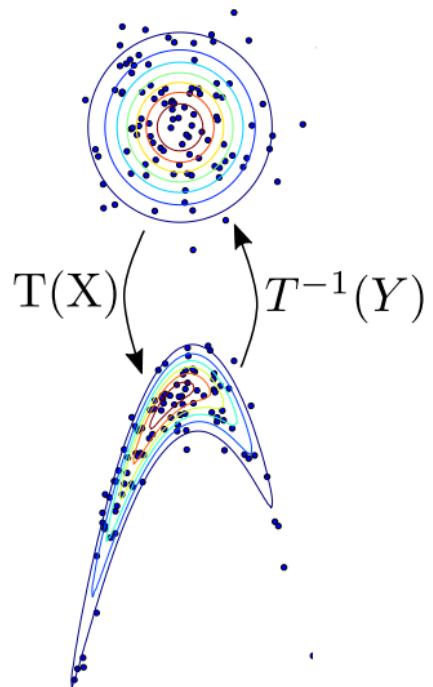
$$\text{PF} \quad T_\sharp \rho = \rho \circ T^{-1} |\nabla T^{-1}|$$

$$\text{PB} \quad T^\sharp \pi = \pi \circ T |\nabla T|$$

- We want T such that

$$\text{PF} \quad \text{For } X \sim \nu_\rho, T(X) \sim \nu_\pi$$

$$\text{PB} \quad \text{For } Y \sim \nu_\pi, T^{-1}(Y) \sim \nu_\rho$$



Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

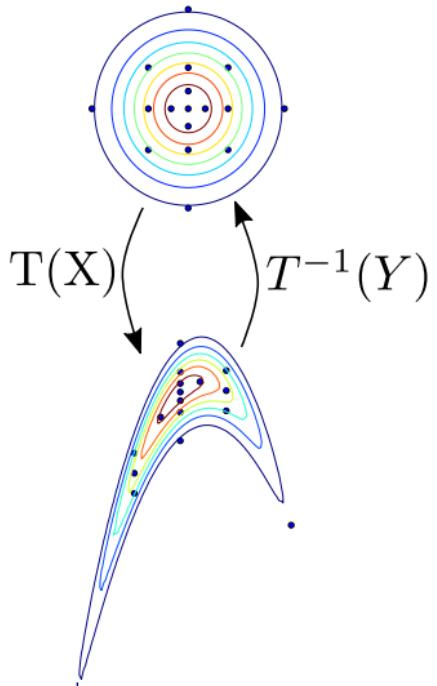
$$\text{PF} \quad T_{\sharp}\rho = \rho \circ T^{-1} |\nabla T^{-1}|$$

$$\text{PB} \quad T^{\sharp}\pi = \pi \circ T |\nabla T|$$

- We want T such that

$$\text{PF} \quad \text{For } X \sim \nu_\rho, T(X) \sim \nu_\pi$$

$$\text{PB} \quad \text{For } Y \sim \nu_\pi, T^{-1}(Y) \sim \nu_\rho$$



Transport maps – Pullbacks [PB] and Pushforwards [PF]

- Distribution ν_ρ with density $\rho : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- Distribution ν_π with density $\pi : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- For $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ we define

$$\text{PF} \quad T_\sharp \rho = \rho \circ T^{-1} |\nabla T^{-1}|$$

$$\text{PB} \quad T^\sharp \pi = \pi \circ T |\nabla T|$$

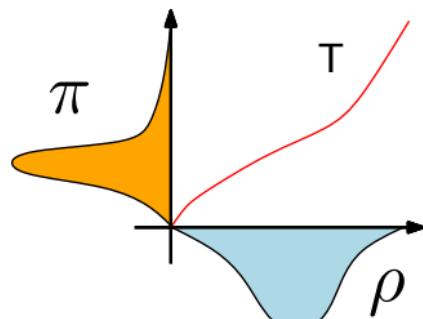
- We want T such that

$$\text{PF} \quad \text{For } X \sim \nu_\rho, T(X) \sim \nu_\pi$$

$$\text{PB} \quad \text{For } Y \sim \nu_\pi, T^{-1}(Y) \sim \nu_\rho$$

Knothe-Rosenblatt rearrangement

$\forall \nu_\rho, \nu_\pi$ Lebesgue absolutely continuous
 \exists a **triangular monotone** map s.t. $T_\sharp \rho = \pi$



$$T(\mathbf{x}) = \begin{bmatrix} T^{(1)}(x_1) \\ T^{(2)}(x_1, x_2) \\ \vdots \\ T^{(d)}(x_1, \dots, x_d) \end{bmatrix}$$

Triangular monotone maps

$$\mathcal{T}_> = \left\{ T : \mathbb{R}^d \rightarrow \mathbb{R}^d : \overbrace{[T(\mathbf{x})]_k = T^{(k)}(x_1, \dots, x_k)}^{\text{triangular}} \text{ and } \overbrace{\partial_{x_k} T^{(k)} > 0}^{\text{monotone}} \right\}$$

Triangular monotone maps

$$\mathcal{T}_> = \left\{ T : \mathbb{R}^d \rightarrow \mathbb{R}^d : \overbrace{[T(\mathbf{x})]_k = T^{(k)}(x_1, \dots, x_k)}^{\text{triangular}} \text{ and } \overbrace{\partial_{x_k} T^{(k)} > 0}^{\text{monotone}} \right\}$$

Integrated squared representation – $\varepsilon > 0$

$$T^{(k)}(x_{1:k}) = c_k(x_{1:k-1}) + \int_0^{x_k} \left(h_k(x_{1:k-1}, t) \right)^2 + \varepsilon dt$$

Triangular monotone maps

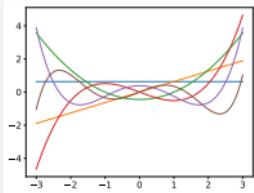
$$\mathcal{T}_>^n = \left\{ T : \mathbb{R}^d \rightarrow \mathbb{R}^d : \overbrace{[T(\mathbf{x})]_k = T^{(k)}(x_1, \dots, x_k)}^{\text{triangular}} \text{ and } \overbrace{\partial_{x_k} T^{(k)} > 0}^{\text{monotone}} \right\}$$

Integrated squared representation – $\varepsilon > 0$

$$T^{(k)}(x_{1:k}) = c_k(x_{1:k-1}) + \int_0^{x_k} \left(h_k(x_{1:k-1}, t) \right)^2 + \varepsilon dt$$

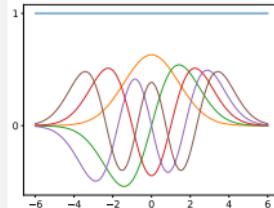
Constant part

$$c_k(x_{1:k-1}) = \sum_{\mathbf{i} \in \mathcal{I}_k} \mathbf{a}_{\mathbf{i}} \Phi_{\mathbf{i}}(x_{1:k-1})$$



Squared part

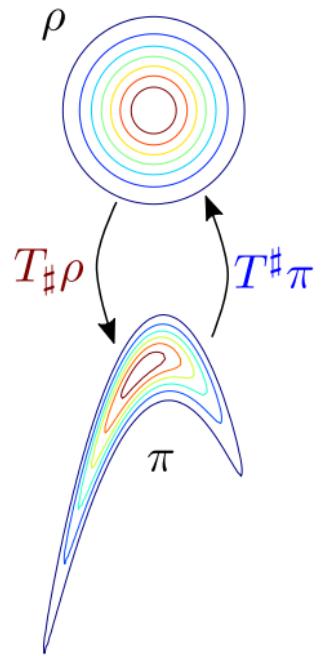
$$h_k(x_{1:k-1}, t) = \sum_{\mathbf{j} \in \mathcal{J}_k} \mathbf{b}_{\mathbf{j}} \Psi_{\mathbf{j}}(x_{1:k-1}, t)$$



Knothe-Rosenblatt rearrangement

$\forall \nu_\rho, \nu_\pi$ Lebesgue absolutely continuous
 \exists a **triangular monotone** map s.t. $T_\# \rho = \pi$

How to find the map $T \in \mathcal{T}_>$
such that $T_\# \rho = \pi$?



Minimize KL-divergence to find optimal map

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_{\sharp} \nu_{\rho} \| \nu_{\pi}) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T_{\sharp} \pi} \right]$$

Minimize KL-divergence to find optimal map

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_\sharp \nu_\rho \| \nu_\pi) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_\rho \left[\log \frac{\rho}{T^\sharp \pi} \right]$$

- + Gradient-based unconstrained optimization if gradients are available
- + We can explore π in parallel
- + We can generate i.i.d. samples from $T^*_\sharp \nu_\rho = \nu_\pi$ in parallel

Minimize KL-divergence to find optimal map

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_\sharp \nu_\rho \| \nu_\pi) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_\rho \left[\log \frac{\rho}{T^\sharp \pi} \right]$$

- + Gradient-based unconstrained optimization if gradients are available
- + We can explore π in parallel
- + We can generate i.i.d. samples from $T^*_\sharp \nu_\rho = \nu_\pi$ in parallel

We are working on $\mathcal{T}_>^n \subset \mathcal{T}_>$, so
how can we evaluate the quality of the approximation?

Convergence criterion – Variance diagnostic

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_{\sharp} \nu_{\rho} \| \nu_{\pi}) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \tilde{\pi}} \right] + \log \int \tilde{\pi}$$

Optimal $T^* \in \mathcal{T}_>$ and $\int \tilde{\pi} = 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{(T^*)^{\sharp} \tilde{\pi}} \right] = 0$

But, optimal $\tilde{T}^* \in \mathcal{T}_>^n$ or $\int \tilde{\pi} \neq 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{(\tilde{T}^*)^{\sharp} \tilde{\pi}} \right] \neq 0$

Convergence criterion – Variance diagnostic

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_{\sharp} \nu_{\rho} \| \nu_{\pi}) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T^{\sharp} \tilde{\pi}} \right] + \log \int \tilde{\pi}$$

Optimal $T^* \in \mathcal{T}_>$ and $\int \tilde{\pi} = 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{(T^*)^{\sharp} \tilde{\pi}} \right] = 0$

But, optimal $\tilde{T}^* \in \mathcal{T}_>^n$ or $\int \tilde{\pi} \neq 1 \quad \Rightarrow \quad \mathbb{E}_{\rho} \left[\log \frac{\rho}{(\tilde{T}^*)^{\sharp} \tilde{\pi}} \right] \neq 0$

$$D_{\text{KL}}(T_{\sharp} \nu_{\rho} \| \nu_{\pi}) \approx \frac{1}{2} \mathbb{V} \left[\log \frac{\rho}{T^{\sharp} \tilde{\pi}} \right] \quad \text{as} \quad T \rightarrow T^*$$

Pros & cons

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T \sharp \rho \| \pi) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T \sharp \pi} \right]$$

- + Gradient-based unconstrained optimization if gradients are available
- + We can explore π in parallel
- + We can generate i.i.d. samples from $T^* \nu_\rho = \nu_\pi$ in parallel
- + We can assess convergence!

Pros & cons

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_\sharp \rho \| \pi) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_\rho \left[\log \frac{\rho}{T^\sharp \pi} \right]$$

- + Gradient-based unconstrained optimization if gradients are available
- + We can explore π in parallel
- + We can generate i.i.d. samples from $T_\sharp^\star \nu_\rho = \nu_\pi$ in parallel
- + We can assess convergence!
- + The map can be used as a preconditioner for other unbiased methods

$$\int f(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x} = \int f(\mathbf{x}) \frac{\pi(\mathbf{x})}{T_\sharp \rho(\mathbf{x})} T_\sharp \rho(\mathbf{x}) d\mathbf{x} = \int f \circ T(\mathbf{x}) \frac{T^\sharp \pi(\mathbf{x})}{\rho(\mathbf{x})} \rho(\mathbf{x}) d\mathbf{x}$$

Pros & cons

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T \sharp \rho \| \pi) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_{\rho} \left[\log \frac{\rho}{T \sharp \pi} \right]$$

- + Gradient-based unconstrained optimization if gradients are available
- + We can explore π in parallel
- + We can generate i.i.d. samples from $T_{\sharp}^* \nu_{\rho} = \nu_{\pi}$ in parallel
- + We can assess convergence!
- + The map can be used as a preconditioner for other unbiased methods
- We need to approximate d functions of up to d variables!

$$T(\mathbf{x}) = \begin{bmatrix} T^{(1)}(x_1) \\ T^{(2)}(x_1, x_2) \\ \vdots \\ T^{(d)}(x_1, \dots, x_d) \end{bmatrix}$$

Pros & cons

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T_\sharp \rho \| \pi) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_\rho \left[\log \frac{\rho}{T^\sharp \pi} \right]$$

- + Gradient-based unconstrained optimization if gradients are available
- + We can explore π in parallel
- + We can generate i.i.d. samples from $T_\sharp^\star \nu_\rho = \nu_\pi$ in parallel
- + We can assess convergence!
- + The map can be used as a preconditioner for other unbiased methods
- We need to approximate d functions of up to d variables!

Sources of low-dimensional structure

- Smoothness
- Marginal independence
- Conditional independence [Spantini '18]
- Low-rank structure

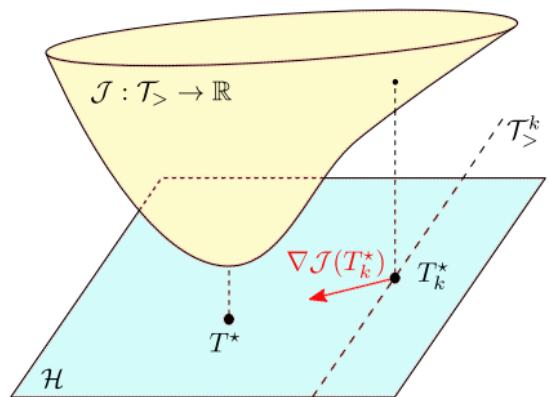
Adaptivity

$$T^* = \arg \min_{T \in \mathcal{T}_>} D_{\text{KL}}(T \sharp \nu_\rho \| \nu_\pi) = \arg \min_{T \in \mathcal{T}_>} \mathbb{E}_\rho \left[\log \frac{\rho}{T \sharp \pi} \right]$$

How to find the **best subset** $\mathcal{T}_>^n \subset \mathcal{T}_>$?

Refinement criterion

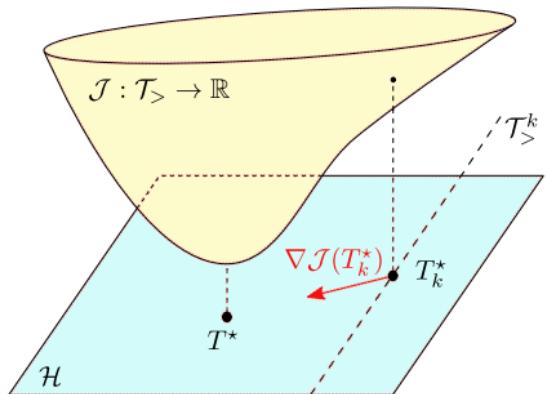
$$\mathbf{a}_k^* = \arg \min_{\mathbf{a} \in \mathbb{R}^{n_k}} \mathbb{E} \left[\underbrace{\log \frac{\rho}{T[\mathbf{a}]^\# \tilde{\pi}}}_{\mathcal{J}(T[\mathbf{a}])} \right]$$



Refinement criterion

$$\mathbf{a}_k^* = \arg \min_{\mathbf{a} \in \mathbb{R}^{n_k}} \mathbb{E} \left[\log \underbrace{\frac{\rho}{T[\mathbf{a}]^\# \tilde{\pi}}}_{\mathcal{J}(T[\mathbf{a}])} \right]$$

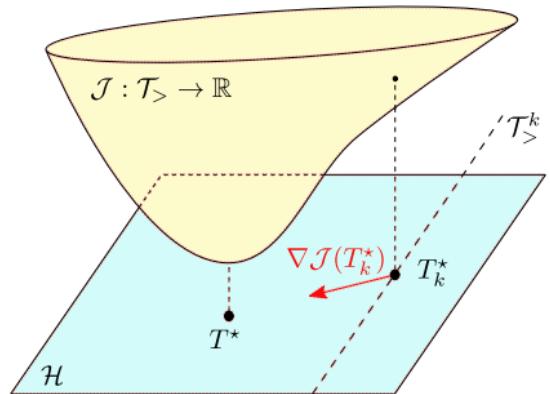
$$\nabla_{\mathbf{a}} \mathcal{J}(T[\mathbf{a}_k^*]) = 0$$



Refinement criterion

$$\mathbf{a}_k^* = \arg \min_{\mathbf{a} \in \mathbb{R}^{n_k}} \mathbb{E} \left[\log \underbrace{\frac{\rho}{T[\mathbf{a}]^\# \tilde{\pi}}}_{\mathcal{J}(T[\mathbf{a}])} \right]$$

$$\nabla_{\mathbf{a}} \mathcal{J}(T[\mathbf{a}_k^*]) = 0$$



The **first variation** $\nabla \mathcal{J}(T[\mathbf{a}_k^*]) \neq 0$

There exists $\varepsilon > 0$ such that

$$\mathcal{J}(T[\mathbf{a}_k^*] - \varepsilon \nabla \mathcal{J}(T[\mathbf{a}_k^*])) < \mathcal{J}(T[\mathbf{a}_k^*])$$

Use the first variation to enrich the approximation space

$$\nabla \mathcal{J}(T[\mathbf{a}_k^*]) = (\nabla_{\mathbf{x}} T)^{-1} \left(\nabla_{\mathbf{x}} \log \frac{\rho}{T[\mathbf{a}_k^*]^\sharp \pi} \right)$$

Use the first variation to enrich the approximation space

$$\nabla \mathcal{J}(T[\mathbf{a}_k^*]) = (\nabla_{\mathbf{x}} T)^{-1} \left(\nabla_{\mathbf{x}} \log \frac{\rho}{T[\mathbf{a}_k^*]^\sharp \pi} \right)$$

- $\nabla \mathcal{J}(T[\mathbf{a}_k^*]) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ **is a map** in $\mathcal{H} \supset \mathcal{T}_>$

Use the first variation to enrich the approximation space

$$\nabla \mathcal{J}(T[\mathbf{a}_k^*]) = (\nabla_{\mathbf{x}} T)^{-1} \left(\nabla_{\mathbf{x}} \log \frac{\rho}{T[\mathbf{a}_k^*]^\sharp \pi} \right)$$

- $\nabla \mathcal{J}(T[\mathbf{a}_k^*]) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a map in $\mathcal{H} \supset \mathcal{T}_>$
- We find the “enriched” space $\mathcal{T}_>^{k+1}$ ($\tilde{\mathcal{T}}^k \supset \mathcal{T}_>^{k+1} \supset \mathcal{T}^k$), using

$$\tilde{\mathbf{b}}_k^* = \arg \min_{\mathbf{b} \in \mathbb{R}^{\tilde{n}_k}} \left\| U[\mathbf{b}] - (T[\mathbf{a}_k^*] - \varepsilon \nabla \mathcal{J}(T[\mathbf{a}_k^*])) \right\|_{L^2_\rho}$$

as a guidance criterion.

Use the first variation to enrich the approximation space

$$\nabla \mathcal{J}(T[\mathbf{a}_k^*]) = (\nabla_{\mathbf{x}} T)^{-1} \left(\nabla_{\mathbf{x}} \log \frac{\rho}{T[\mathbf{a}_k^*]^\sharp \pi} \right)$$

- $\nabla \mathcal{J}(T[\mathbf{a}_k^*]) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a map in $\mathcal{H} \supset \mathcal{T}_>$
- We find the “enriched” space $\mathcal{T}_>^{k+1}$ ($\tilde{\mathcal{T}}^k \supset \mathcal{T}_>^{k+1} \supset \mathcal{T}^k$), using

$$\tilde{\mathbf{b}}_k^* = \arg \min_{\mathbf{b} \in \mathbb{R}^{\tilde{n}_k}} \left\| U[\mathbf{b}] - (T[\mathbf{a}_k^*] - \varepsilon \nabla \mathcal{J}(T[\mathbf{a}_k^*])) \right\|_{L^2_\rho}$$

as a guidance criterion.

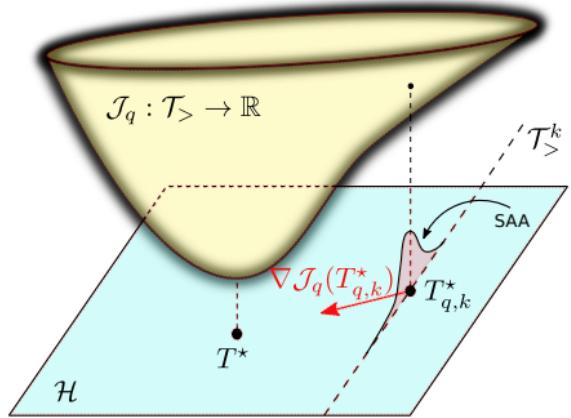
No new evaluation of $\log \pi$ or $\nabla_{\mathbf{x}} \log \pi$ are required

Controlling the sample average accuracy

$$T_k^* = \arg \min_{T \in \mathcal{T}_>^k} -\mathbb{E}_{\rho} [\log T^\sharp \pi] \approx \arg \min_{T \in \mathcal{T}_>^k} \overbrace{-\mathcal{Q}_q [\log T^\sharp \pi(\mathbf{x}_i)]}^{\mathcal{J}_q(T)} =: T_{q,k}^*$$

Controlling the sample average accuracy

$$T_k^* = \arg \min_{T \in \mathcal{T}_>^k} -\mathbb{E}_{\rho} [\log T^\sharp \pi] \approx \arg \min_{T \in \mathcal{T}_>^k} \overbrace{-\mathcal{Q}_q [\log T^\sharp \pi(\mathbf{x}_i)]}^{\mathcal{J}_q(T)} =: T_{q,k}^*$$

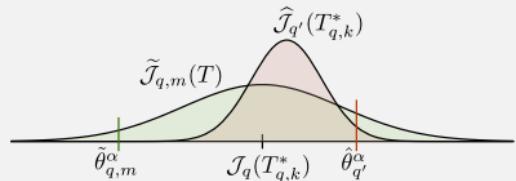


Controlling the sample average accuracy

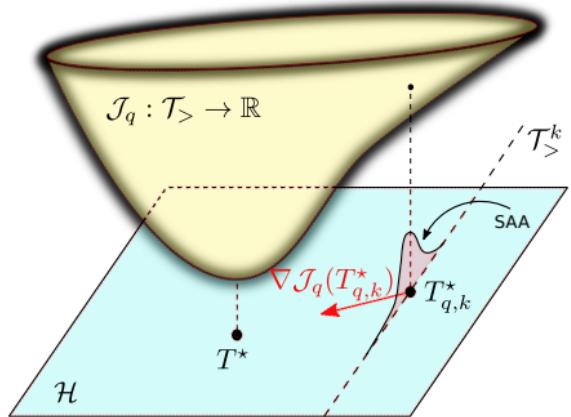
$$T_k^* = \arg \min_{T \in \mathcal{T}_>^k} -\mathbb{E}_{\rho} [\log T^\sharp \pi] \approx \arg \min_{T \in \mathcal{T}_>^k} \overbrace{-\mathcal{Q}_q [\log T^\sharp \pi(\mathbf{x}_i)]}^{\mathcal{J}_q(T)} =: T_{q,k}^*$$

Sample average approximation

$$\tilde{\theta}_{q,m} \leq \mathcal{J}(T_k^*) \leq \hat{\theta}_{q'}$$

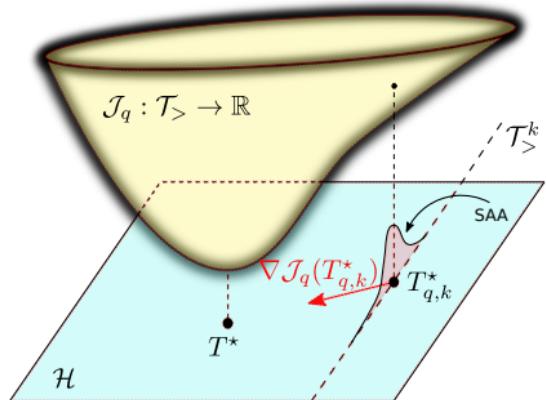
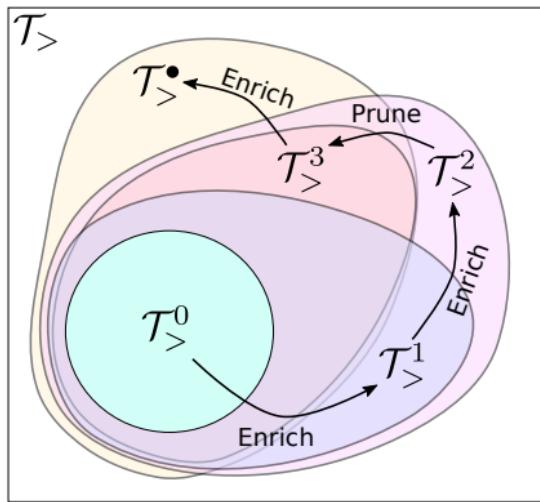


$$\widetilde{J}_{q,m}(T) = \frac{1}{m} \sum_{i=1}^m \min_{T \in \mathcal{T}_>^k} \mathcal{J}_q(T)$$



Adaptivity ingredients

- **Convergence criterion – Variance diagnostic** : $\mathbb{V} \left[\log \frac{\rho}{T^{\sharp} \pi} \right]$
- **Refinement criterion – First variation** : $\nabla \mathcal{J}(T[\mathbf{a}^*])$
- **Stability criterion – Sample average approximation** : $\tilde{\theta}_{q,m} \leq \mathcal{J}(T_k^*) \leq \hat{\theta}_q'$



Off-line learning / On-line inference

Joint distribution $\nu_{\pi'}$: $\pi'(\mathbf{X}, \Theta) := \pi(\mathbf{X}|\Theta) \pi(\Theta)$

Posterior distribution ν_{π} : $\tilde{\pi}(\Theta|\mathbf{X} = \hat{\mathbf{x}}) \propto \mathcal{L}_{\hat{\mathbf{x}}}(\Theta) \pi(\Theta)$

Off-line learning / On-line inference

Joint distribution $\nu_{\pi'}$: $\pi'(\mathbf{X}, \Theta) := \pi(\mathbf{X}|\Theta) \pi(\Theta)$

Posterior distribution ν_{π} : $\tilde{\pi}(\Theta|\mathbf{X} = \hat{\mathbf{x}}) \propto \mathcal{L}_{\hat{\mathbf{x}}}(\Theta) \pi(\Theta)$

Let $T(\mathbf{y}, \theta) = \begin{bmatrix} T^{(\mathbf{y})}(\mathbf{y}) \\ T^{(\theta)}(\mathbf{y}, \theta) \end{bmatrix}$ be such that $T \sharp \mathcal{N}(\mathbf{0}, \mathbf{I}) \approx \nu_{\pi'}$,

then given data $\hat{\mathbf{x}}$

$C_{\hat{\mathbf{y}}}(\theta) := T^{(\theta)}(\hat{\mathbf{y}}, \theta)$ is such that $(C_{\hat{\mathbf{y}}}) \sharp \nu_{\rho} \approx \nu_{\pi}$

where $\hat{\mathbf{y}} = \left(T^{(\mathbf{y})} \right)^{-1}(\hat{\mathbf{x}})$

Biochemical Oxygen Demand



Biochemical Oxygen Demand

We model the oxygen level at time t by

$$X(t) = A(1 - \exp(-Bt)) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2),$$
$$A \sim \mathcal{U}(0.4, 1.2) \quad \text{and} \quad B \sim \mathcal{U}(0.01, 0.31),$$

and we want to

Biochemical Oxygen Demand

We model the oxygen level at time t by

$$X(t) = A(1 - \exp(-Bt)) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2), \\ A \sim \mathcal{U}(0.4, 1.2) \quad \text{and} \quad B \sim \mathcal{U}(0.01, 0.31),$$

and we want to

- ① **Off-line:** characterize the joint distribution $(X(1), \dots, X(4), A, B) \sim \nu_{\pi'}$.

Biochemical Oxygen Demand

We model the oxygen level at time t by

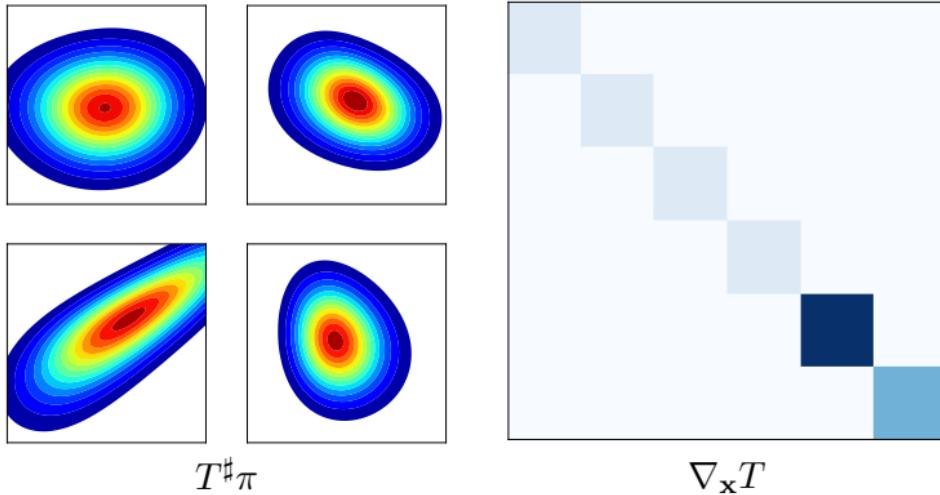
$$X(t) = A(1 - \exp(-Bt)) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2),$$
$$A \sim \mathcal{U}(0.4, 1.2) \quad \text{and} \quad B \sim \mathcal{U}(0.01, 0.31),$$

and we want to

- ① **Off-line:** characterize the joint distribution $(X(1), \dots, X(4), A, B) \sim \nu_{\pi'}$.
- ② **On-line:** for data $\hat{\mathbf{x}}$, generate samples from the posterior $(A, B | \mathbf{X} = \hat{\mathbf{x}}) \sim \nu_{\pi}$.

Biochemical Oxygen Demand – off-line phase

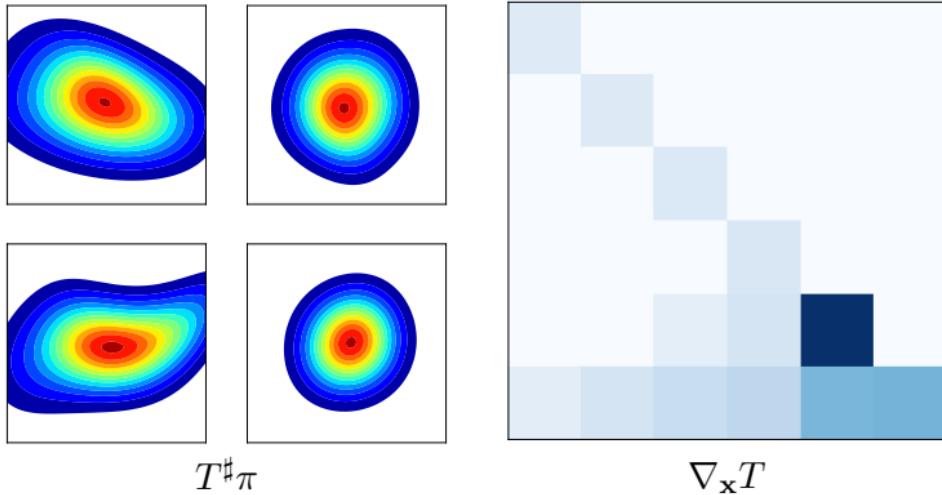
Iteration 1



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

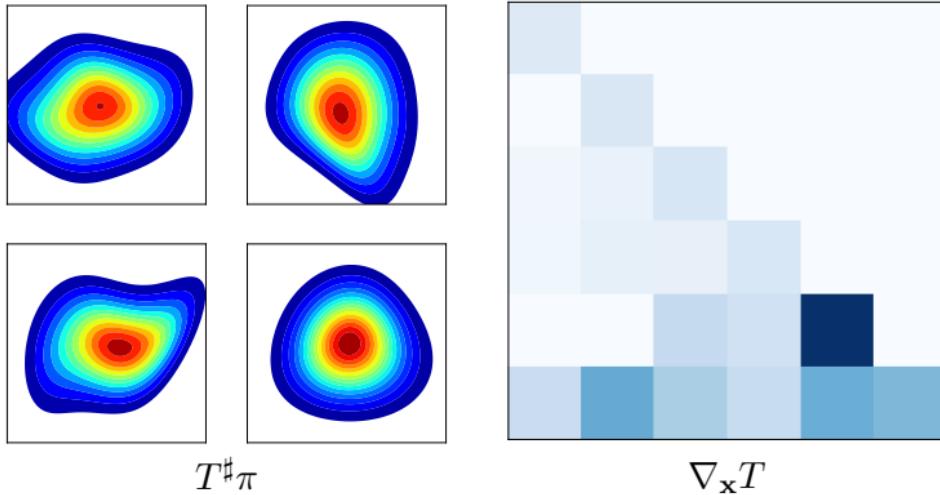
Iteration 2



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

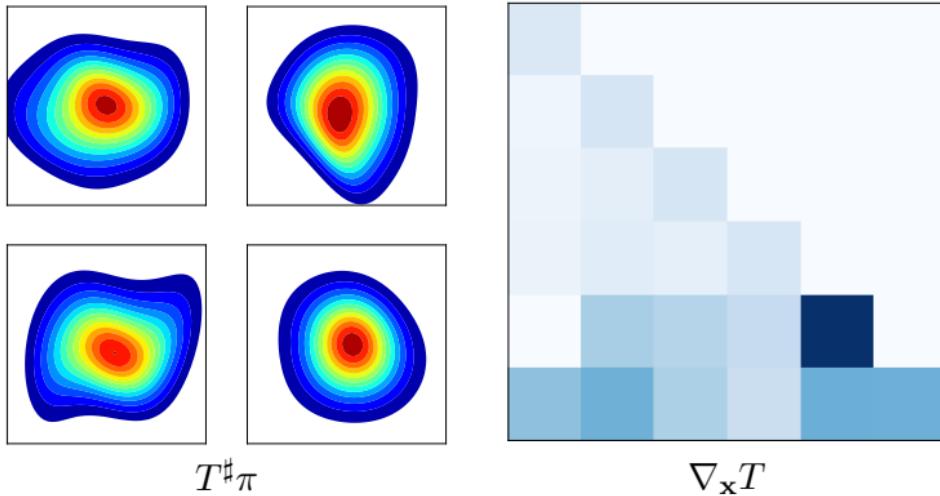
Iteration 3



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

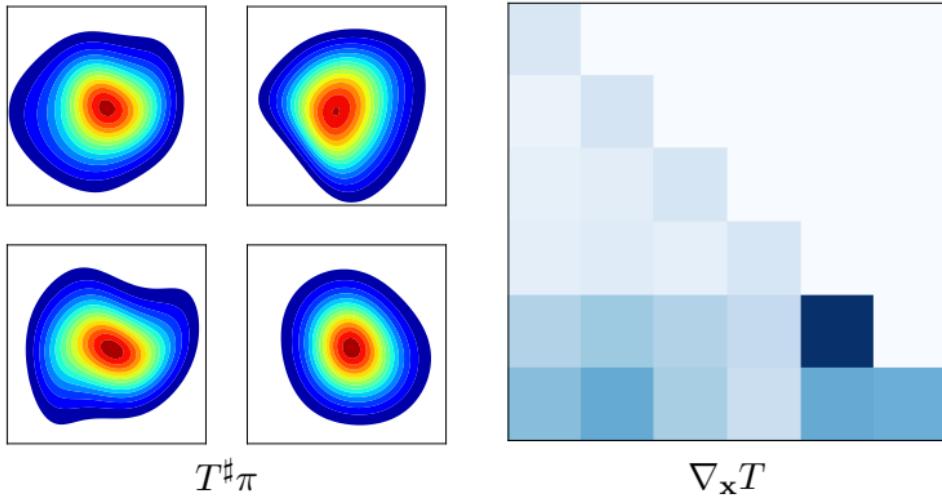
Iteration 4



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

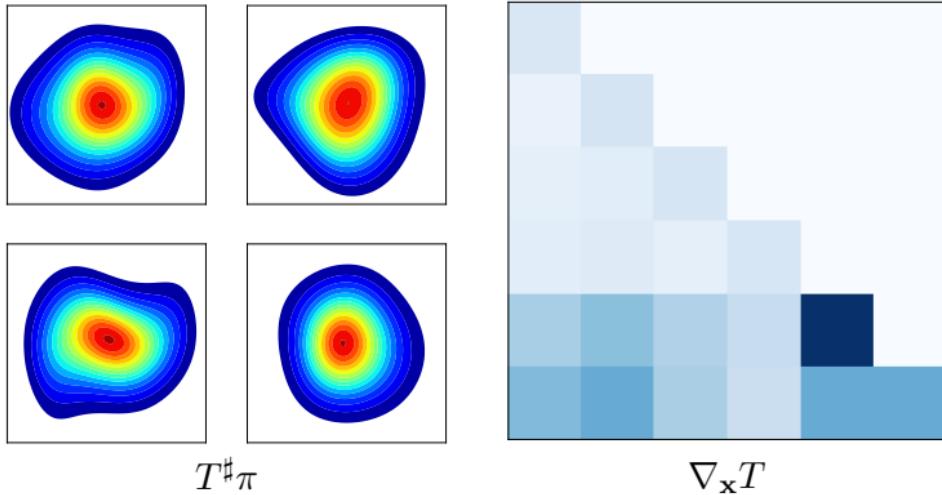
Iteration 5



Reminder: $T^\sharp\pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

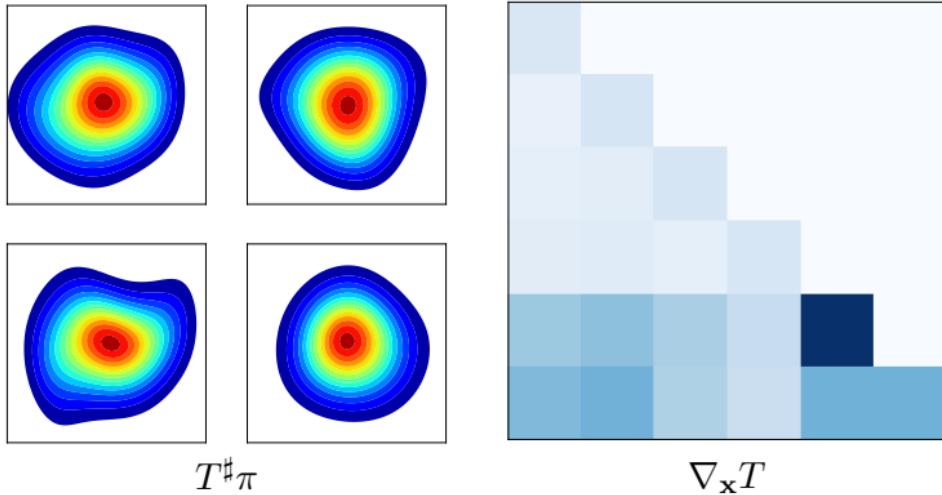
Iteration 6



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

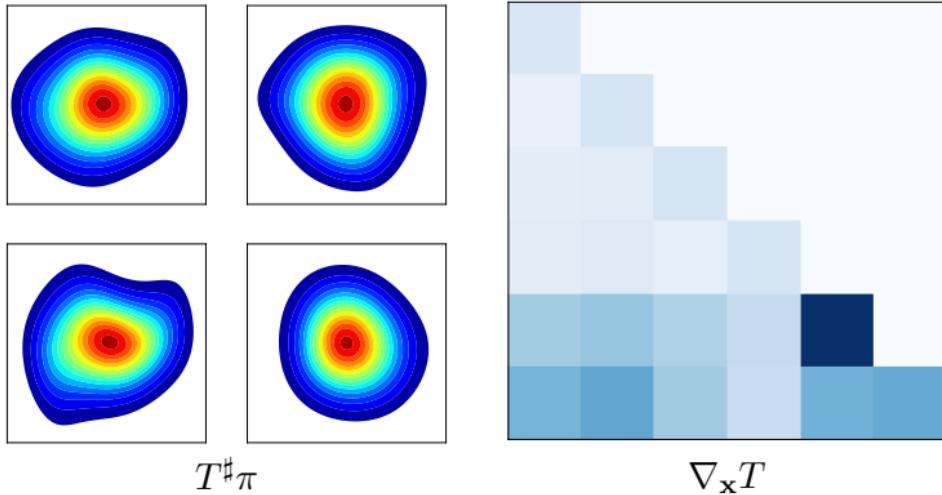
Iteration 7



Reminder: $T^{\sharp}\pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

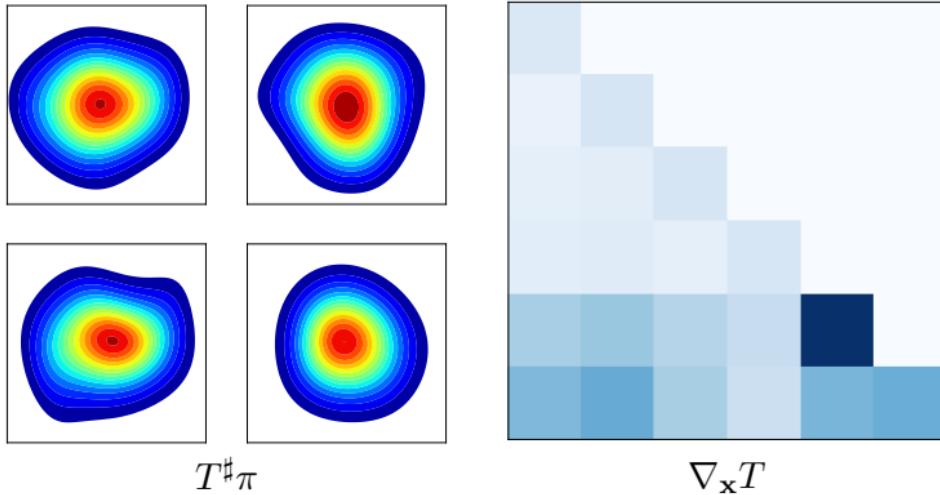
Iteration 8



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

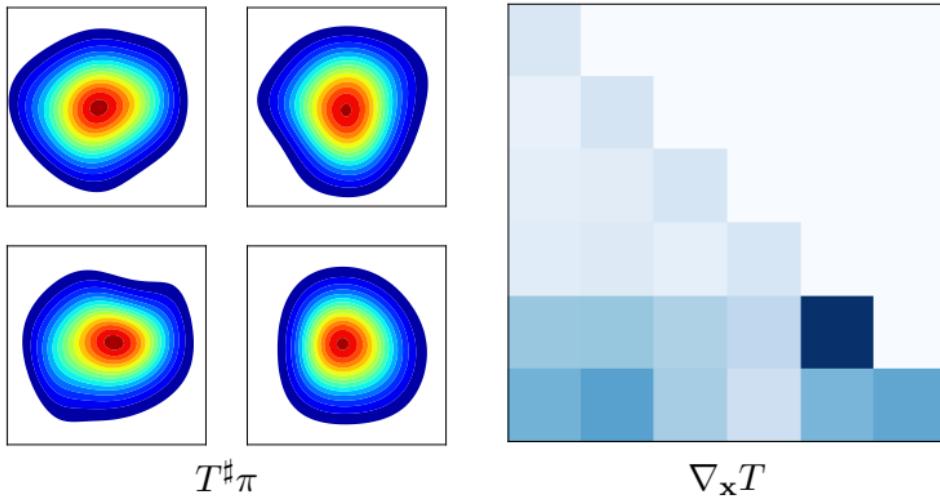
Iteration 9



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

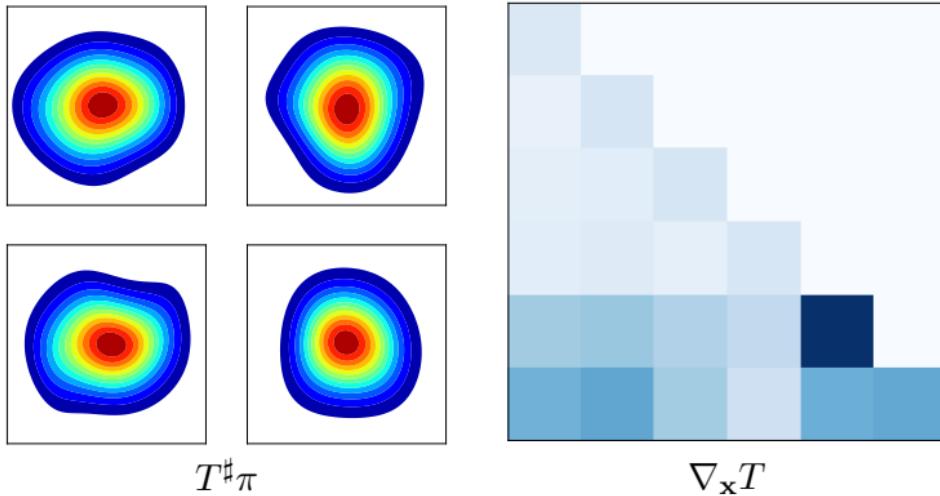
Iteration 10



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

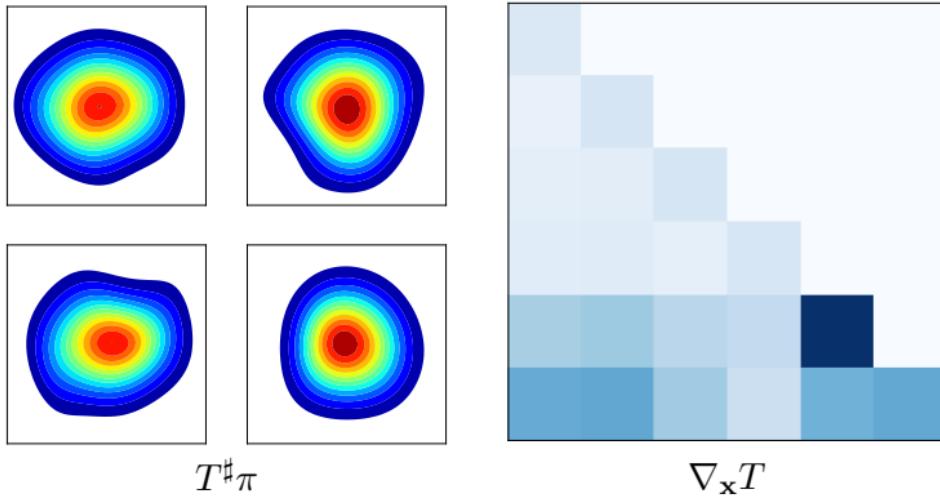
Iteration 11



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

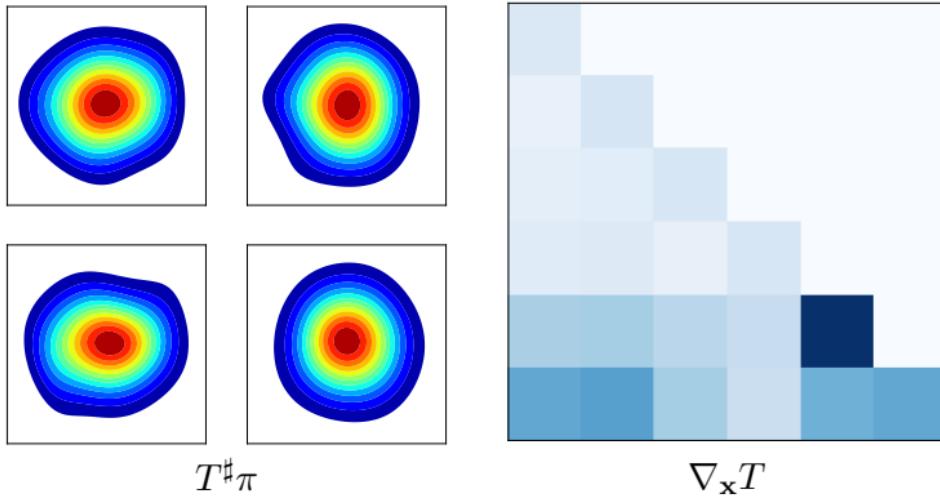
Iteration 12



Reminder: $T^\sharp \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

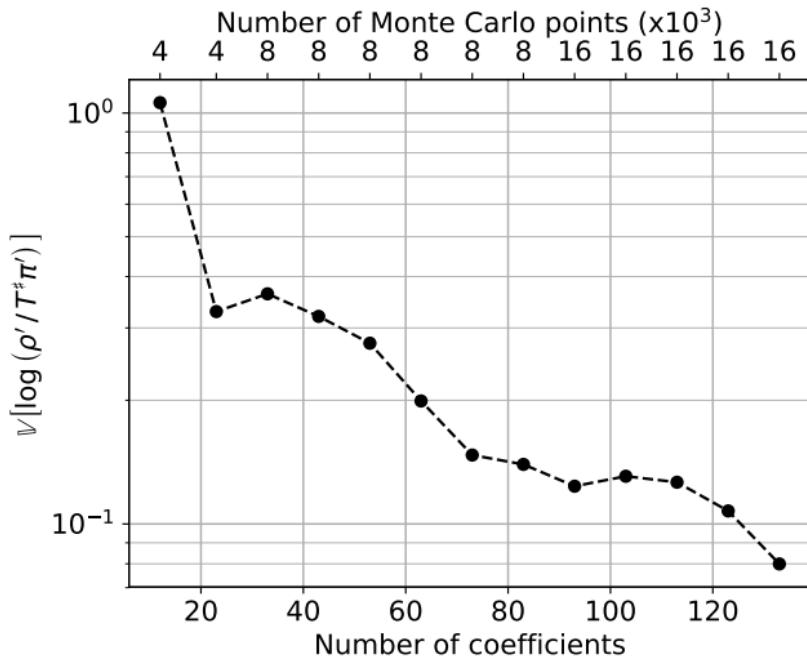
Biochemical Oxygen Demand – off-line phase

Iteration 13

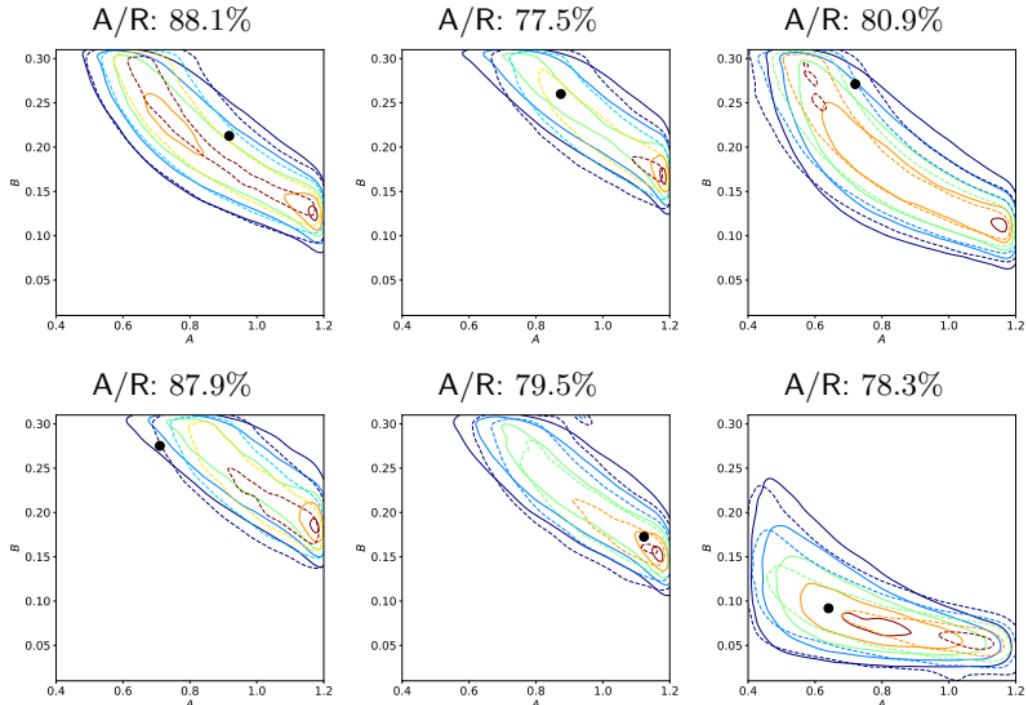


Reminder: $T^\# \pi' \approx \rho'$, where ρ' is the density of $\mathcal{N}(0, \mathbf{I})$

Biochemical Oxygen Demand – off-line phase

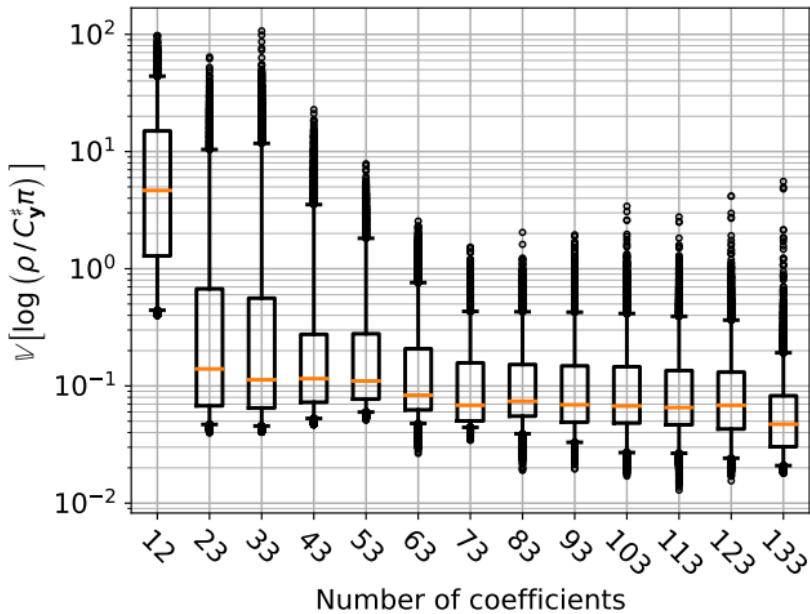


Biochemical Oxygen Demand – on-line phase



Dashed lines are contours of Markov Chains long 10^6
generated with Metropolis-Hastings with independent proposals.

Biochemical Oxygen Demand – on-line phase



Key contributions

Algorithms for characterizing probability measures
via **deterministic couplings** and **optimization**,
exploiting **smoothness** and **marginal independence**

Contact: Daniele Bigoni – **dabi@mit.edu**

Software: <https://transportmaps.mit.edu>

Bigoni et al. “Adaptive construction of measure transports for Bayesian inference”
Spantini et al. “Inference via low-dimensional couplings”

References: Marzouk et al. “An introduction to sampling via measure transport”
Parno et al. “Transport map accelerated Markov chain Monte Carlo”
El Moselhy et al. “Bayesian inference with optimal maps”

Thanks to:

