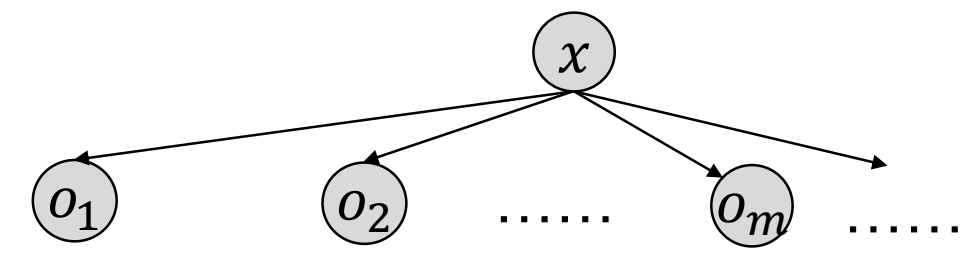


## Bayes' Rule

Given

- 1 **Prior** distribution  $\pi(\mathbf{x})$
- 2 **Likelihood** function  $p(\mathbf{o}|\mathbf{x})$
- 3 **Observations**  $\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_m$



The **posterior distribution** of unknown parameter  $\mathbf{x}$  can be computed by Bayes' Rule:

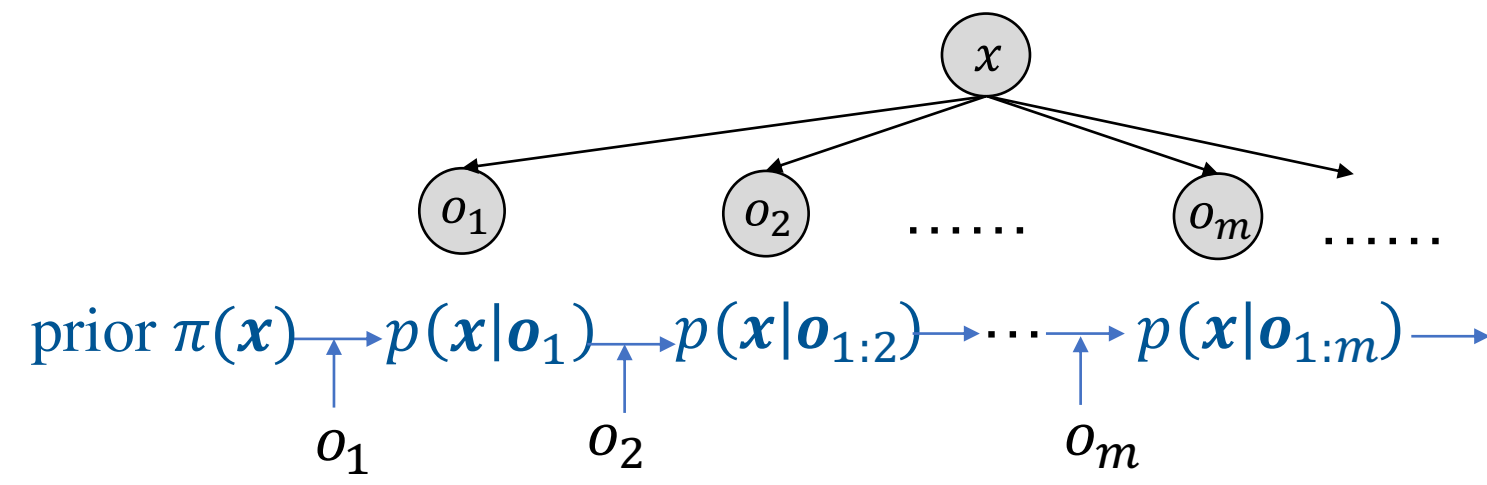
$$p(\mathbf{x}|\mathbf{o}_{1:m}) = \frac{1}{z} \pi(\mathbf{x}) \prod_{i=1}^m p(\mathbf{o}_i|\mathbf{x})$$

$$z = \int \pi(\mathbf{x}) \prod_{i=1}^m p(\mathbf{o}_i|\mathbf{x}) d\mathbf{x}$$

Challenging computational problem for high dimensional  $\mathbf{x}$

## Sequential Bayesian Inference

Observations  $\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_m$  arrive sequentially:



An ideal algorithm should:

- **Efficiently update**  $p(\mathbf{x}|\mathbf{o}_{1:m})$  to  $p(\mathbf{x}|\mathbf{o}_{1:m+1})$  when  $\mathbf{o}_{m+1}$  is observed
- **Without storing** all historical observations  $\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_m$

$$p(\mathbf{x}|\mathbf{o}_{1:m}) \propto p(\mathbf{x}|\mathbf{o}_{1:m-1}) p(\mathbf{o}_m|\mathbf{x})$$

updated posterior      current posterior      likelihood

## Related Works

MCMC

- requires a complete scan of the data

Variational Inference (VI)

- requires re-optimization for every new observation

Stochastic approximate inference

- are prescribed algorithms to optimize the final posterior  $p(\mathbf{x}|\mathbf{o}_{1:M})$
- can not exploit the structure of the sequential inference problem

Sequential monte Carlo

- state of art for online Bayesian Inference
- but suffers from path degeneracy problem in high dimensions
- rejuvenation steps can help but will violate online constraints

An Operator View: Kernel Bayes' Rule

- the posterior is represented as an embedding  $\mu_m = \mathbb{E}_{p(\mathbf{x}|\mathbf{o}_{1:m})} \phi(\mathbf{x})$

$$\mu_{m+1} = \mathcal{K}(\mu_m, \mathbf{o}_{m+1})$$

updated embedding      current embedding

- views the Bayes update as an operator in RKHS

## Our Method

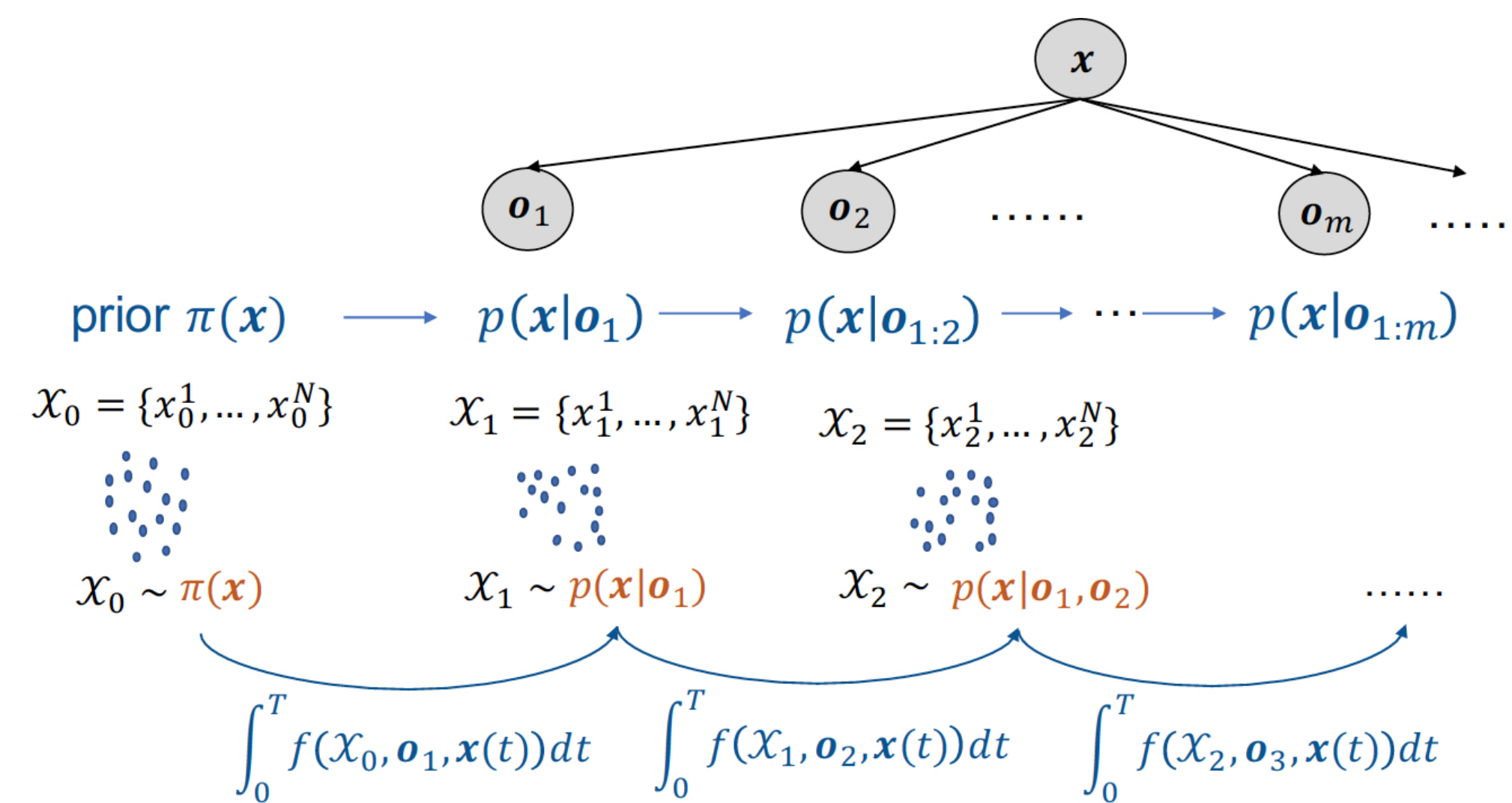
Start with  $N$  particles

$$\mathcal{X}_0 = \{x_0^1, x_0^2, \dots, x_0^N\}, \text{ sampled i.i.d. from prior } \pi(\mathbf{x})$$

Transport particles to next posterior as the solution of ODEs

$$\begin{cases} \frac{d\mathbf{x}}{dt} = f(\mathcal{X}_0, \mathbf{o}_1, \mathbf{x}(t), t), \forall t \in (0, T] \\ \mathbf{x}(0) = \mathbf{x}_0^n \end{cases} \xrightarrow{\text{gives}} \mathbf{x}_1^n = \mathbf{x}(T)$$

## Particle Flow Bayes' Rule



Particle Flow as a Bayesian Operator

$$\mathbf{x}_{m+1}^n = \mathcal{F}(\mathcal{X}_m, \mathbf{o}_{m+1}, \mathbf{x}_m^n) := \mathbf{x}_m^n + \int_0^T f(\mathcal{X}_m, \mathbf{o}_{m+1}, \mathbf{x}(t), t) dt.$$

$$\log q_{m+1}(\mathbf{x}_{m+1}^n) = \log q_m(\mathbf{x}_m^n) - \int_0^T \nabla_{\mathbf{x}} \cdot f dt.$$

## Advantages: Flow Property

There are mainly two obvious advantages of Particle Flow:

- 1 First, the **location** of the particles can be **moved** according to posterior distribution.
- 2 Second, the probability density can be computed efficiently because the **change of log-density also follows a ODE**.

- **Continuity Equation** express the law of *local conservation of mass*: (1) Mass can neither be created nor destroyed; (2) nor can it 'teleport' from one place to another.

$$\frac{\partial q(\mathbf{x}, t)}{\partial t} = -\nabla_{\mathbf{x}} \cdot (qf)$$

- **Theorem.** If  $\frac{d\mathbf{x}}{dt} = f$ , then the change in log-density follows

$$\frac{d \log q(\mathbf{x}, t)}{dt} = -\nabla_{\mathbf{x}} \cdot f.$$

## Does A Unified Flow Velocity $f$ exist?

$$\mathbf{x}(0) \sim \pi(\mathbf{x}) \quad \mathbf{x}(t) \sim p(\mathbf{x}|\mathbf{o}_1)$$

$$\mathbf{x}(T) = \mathbf{x}(0) + \int_0^T f(\text{inputs}) dt$$

Does a **unified flow velocity**  $f$  exist for different Bayesian inference tasks involving different priors and different observations?

## Existence of Flow-based Bayes' Rule

(1) **Langevin dynamics** is a *stochastic* process

$$d\mathbf{x}(t) = \nabla_{\mathbf{x}} \log \pi(\mathbf{x}) p(\mathbf{o}|\mathbf{x}) dt + \sqrt{2} d\mathbf{w}(t),$$

where  $d\mathbf{w}(t)$  is a standard Brownian motion.

- The probability density  $q(\mathbf{x}, t)$  of  $\mathbf{x}(t)$  converges to a stationary distribution, which is the posterior  $p(\mathbf{x}|\mathbf{o})$ .

(2) **Stochastic Flow to Deterministic Flow.**

- The probability density  $q(\mathbf{x}, t)$  of Langevin dynamics follows a **deterministic** evolution according to **Fokker-Planck equation**

$$\frac{\partial q}{\partial t} = -\nabla_{\mathbf{x}} \cdot (q \nabla_{\mathbf{x}} \log \pi(\mathbf{x}) p(\mathbf{o}|\mathbf{x})) + \nabla_{\mathbf{x}} q(\mathbf{x}, t).$$

- Fokker-Planck equation can be rewritten in the form of **Continuity Equation**:

$$\frac{\partial q}{\partial t} = -\nabla_{\mathbf{x}} \cdot (qf),$$

where  $f = \nabla_{\mathbf{x}} \log \pi(\mathbf{x}) p(\mathbf{o}|\mathbf{x}) - \nabla_{\mathbf{x}} \log q(\mathbf{x}, t)$ .

$\Rightarrow$  **deterministic flow!**

(3) **Closed-Loop to Open-Loop:** The above deterministic flow is closed-loop, which depends on flow state  $q(\mathbf{x}, t)$ . We use optimal control theory to show there exists a unified  $f$  which is independent of  $q(\mathbf{x}, t)$ .

**Conclusion of a unified  $f$ .** There exists a fixed and deterministic flow velocity  $f$  of the form

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{o}_{1:m}) p(\mathbf{o}_{m+1}|\mathbf{x}) - w^*(p(\mathbf{x}|\mathbf{o}_{1:m}), t),$$

which can transform  $p(\mathbf{x}|\mathbf{o}_{1:m})$  to  $p(\mathbf{x}|\mathbf{o}_{1:m+1})$  and in turns define a unified particle flow Bayes operator  $\mathcal{F}$ .

## Parameterization

$$f(p(\mathbf{x}|\mathbf{o}_{1:m}), p(\mathbf{o}_{m+1}|\mathbf{x}), \mathbf{x}(t), t) \Rightarrow f(\mathcal{X}_m, \mathbf{o}_{m+1}, \mathbf{x}(t), t)$$

- $p(\mathbf{x}|\mathbf{o}_{1:m}) \Rightarrow \mathcal{X}_m$   
Use samples  $\mathcal{X}_m$  as surrogates, feature space embedding.
- $p(\mathbf{o}_{m+1}|\mathbf{x}) \Rightarrow (\mathbf{o}_{m+1}, \mathbf{x}(t))$

Overall we parameterize the flow velocity as

$$f = \mathbf{h} \left( \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_m^n), \mathbf{o}_{m+1}, \mathbf{x}(t), t \right),$$

where  $\mathbf{h}$  and  $\phi$  are neural networks. Let  $\theta \in \Theta$  be their parameters which are independent of  $t$ .

## Learning Algorithm

**Multi-task Framework**

- The training set  $\mathcal{D}_{\text{train}}$  contains **multiple inference tasks**
- Each task  $\mathcal{T} \in \mathcal{D}_{\text{train}}$  is a tuple

$$\mathcal{T} := (\underbrace{\pi(\mathbf{x})}_{\text{prior}}, \underbrace{p(\cdot|\mathbf{x})}_{\text{likelihood}}, \underbrace{\{\mathbf{o}_1, \dots, \mathbf{o}_M\}}_M \text{ observations})$$

**Loss Function**

- The loss for each  $\mathcal{T}$  is  $\sum_{m=1}^M \text{KL}(q_m(\mathbf{x}) || p(\mathbf{x}, \mathbf{o}_{1:m}))$ , where  $q_m(\mathbf{x})$  is the distribution transported by  $\mathcal{F}$  at  $m$ -th stage.

- Equivalent to minimize negative evidence lower bound (ELBO)

$$\mathcal{L}(\mathcal{T}) = \sum_{m=1}^M \sum_{n=1}^N (\log q_m(\mathbf{x}_m^n) - \log p(\mathbf{x}_m^n, \mathbf{o}_{1:m})).$$

- Cumulative loss:  $\mathcal{L}(\mathcal{D}_{\text{train}}) = \sum_{\mathcal{T} \in \mathcal{D}_{\text{train}}} \mathcal{L}(\mathcal{T})$ .

## Experiment 1: Benefits for High Dimension

Multivariate Gaussian Model

- prior  $\mathbf{x} \sim \mathcal{N}(\mu_x, \Sigma_x)$
- observation conditioned on prior  $\mathbf{o}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, \Sigma_o)$

Experiment Setting

- Training set only contains *sequences of 10 observations*.
- Testing set contains 25 difference *sequences of 100 observations*.

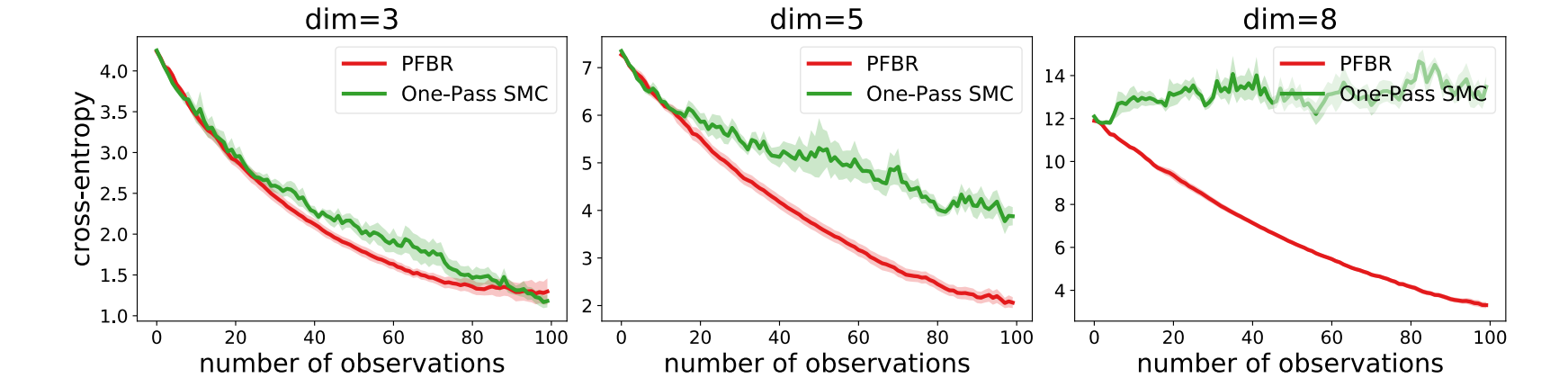


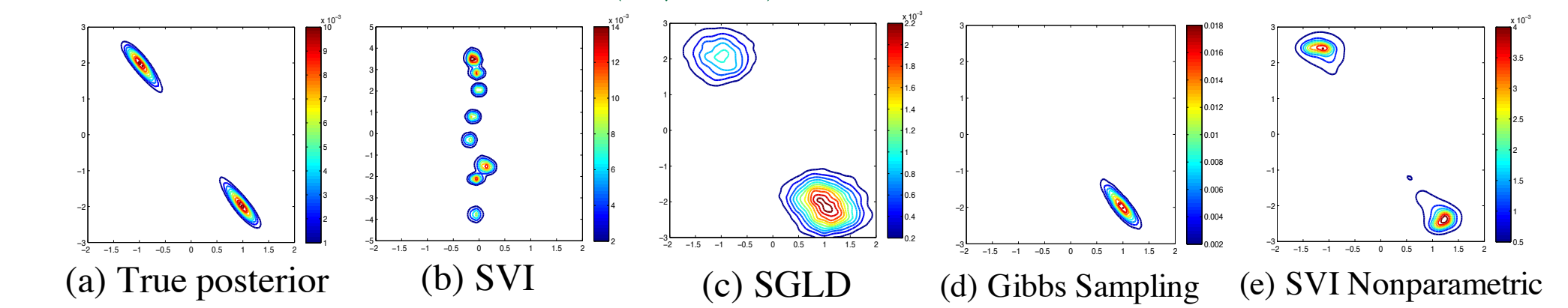
Figure: Cross entropy  $\mathbb{E}_{p(\mathbf{x}|\mathbf{o}_{1:m})} - \log q_m$

## Experiment 2: Multi-Modal Posterior

Gaussian Mixture Model

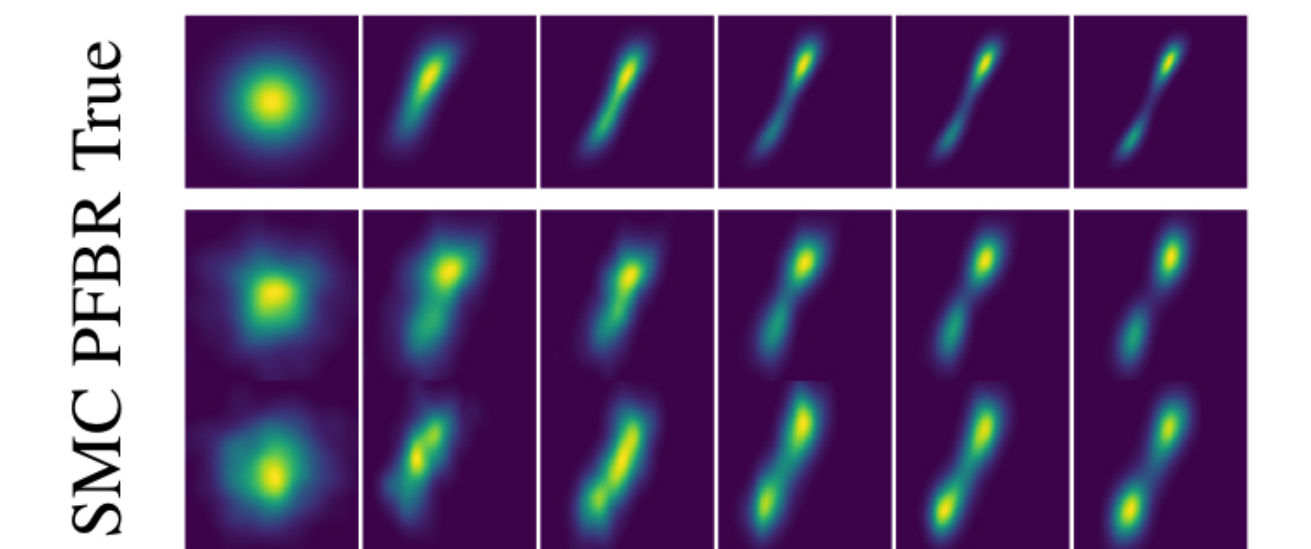
- prior  $\mathbf{x}_1, \mathbf{x}_2 \sim \mathcal{N}(0, 1)$
- observations  $\mathbf{o}|\mathbf{x}_1, \mathbf{x}_2 \sim \frac{1}{2} \mathcal{N}(\mathbf{x}_1, 1) + \frac{1}{2} \mathcal{N}(\mathbf{x}_1 + \mathbf{x}_2, 1)$
- With  $(\mathbf{x}_1, \mathbf{x}_2) = (1, -2)$ , the posterior has two modes.

To fit only one posterior  $p(\mathbf{x}|\mathbf{o}_{1:m})$  is already not easy.



Our more challenging experimental setting:

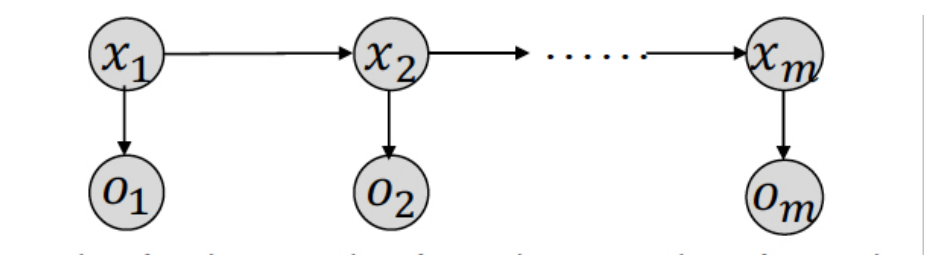
- The learned MPF operator will be tested on sequences that are not observed in training set
- It needs to fit all *intermediate posteriors*  $p(\mathbf{x}|\mathbf{o}_1), p(\mathbf{x}|\mathbf{o}_1, \mathbf{o}_2), \dots$



## Experiment 3: Hidden Markov Model

Hidden Markov Model - Linear Dynamical System

- $\mathbf{x}_m = A\mathbf{x}_{m-1} + \epsilon_m, \epsilon_m \sim \mathcal{N}(0, \Sigma_1)$
- $\mathbf{o}_m = B\mathbf{x}_m + \delta_m, \delta_m \sim \mathcal{N}(0, \Sigma_2)$



Marginal posteriors update:  $\rightarrow p(x_1|\mathbf{o}_1) \rightarrow p(x_2|\mathbf{o}_{1:2}) \rightarrow p(x_m|\mathbf{o}_{1:m})$

Transition sampling + PFBR operator:

$$(i) \tilde{\mathbf{x}}_m^n = A\mathbf{x}_{m-1}^n + \epsilon_m, (ii) \mathbf{x}_m^n = \mathcal{F}(\tilde{\mathcal{X}}_m^n, \tilde{\mathbf{x}}_m^n, \mathbf{o}_{m+1}),$$

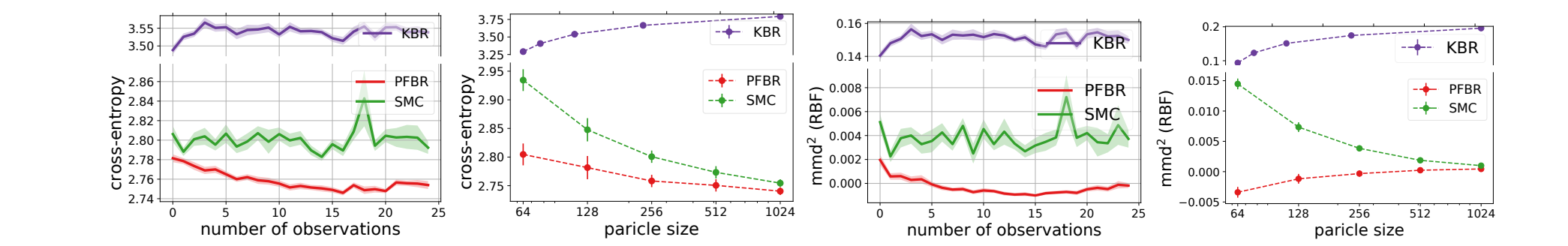


Figure: Left: cross entropy; Right: Square MMD with RBK kernel.

## Experiment 4: Bayesian Logistic Regression

BLR on MNIST dataset 8 vs 6

- Likelihood function  $p(o_m|\theta) = y^{o_m} (1-y)^{1-o_m}$ , where  $y = \sigma(\theta^T \phi_m)$ .

Multi-task Environment

- Reduce the dimension 50 by PCA
- Rotate the first two components by an angle  $\psi \sim [-15^\circ, 15^\circ]$

predict  $\rightarrow$  observe true labels  $\rightarrow$  update particles  $\rightarrow$  predict  $\rightarrow$  observe true labels  $\rightarrow \dots$

