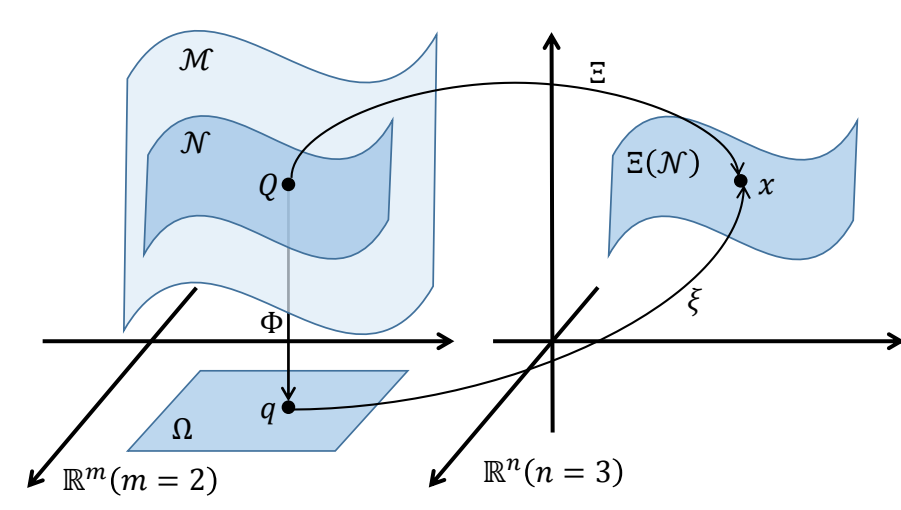


# STOCHASTIC GRADIENT GEODESIC MCMC METHODS



CHANG LIU, JUN ZHU, YANG SONG

{chang-li14@mails., dcszj@}tsinghua.edu.cn; songyang@stanford.edu



## INTRODUCTION

### Task

Scalable Bayesian inference for latent variables on Riemann manifolds by Monte Carlo.

### Drawbacks of Existing Methods

- Unscalable: drawing one sample needs traversing the whole dataset.
- Inner iteration: iteration within one dynamics simulation step.
- Global coordinates requirement: limited applicability (fail for e.g. hypersphere).
- Lower order integrator.

### Our Solution

Stochastic Gradient Geodesic Monte Carlo; geodesic SG Nosé-Hoover Thermostats.

Table: A summary of related methods (-: not for manifold r.v.; †: not SSI; ‡: 2nd-order versions appear afterwards)

methods	stochastic gradient	no inner iteration	no global coordinates	order of integrator
GMC	×	✓	✓	2nd
RMLD	×	✓	×	1st
RMHMC	×	×	×	2nd <sup>†</sup>
CHMC	×	×	✓	2nd <sup>†</sup>
SGLD	✓	✓	-	1st
SGHMC	✓	✓	-	1st <sup>‡</sup>
SGNHT	✓	✓	-	1st <sup>‡</sup>
SGRLD	✓	✓	×	1st
SGRHMC	✓	✓	×	1st
SGGMC	✓	✓	✓	2nd
gSGNHT	✓	✓	✓	2nd

## PRELIMINARIES

### SG-MCMC

To sample from posterior  $\pi(q|\mathcal{D})$ , estimate the required gradient  $\nabla U(q) \triangleq -\nabla \log \pi(q|\mathcal{D})$

$$= -\nabla \log \pi_0(q) - \sum_{d=1}^D \nabla \log \pi(x_d|q)$$

by stochastic gradient with random subset  $\mathcal{S}$ :

$$\nabla \tilde{U}(q) \triangleq -\nabla \log \pi_0(q) - (D/|\mathcal{S}|) \sum_{x \in \mathcal{S}} \nabla \log \pi(x|q).$$

- A complete recipe (Ma et al., 2015) for the dynamics of SG-MCMC:

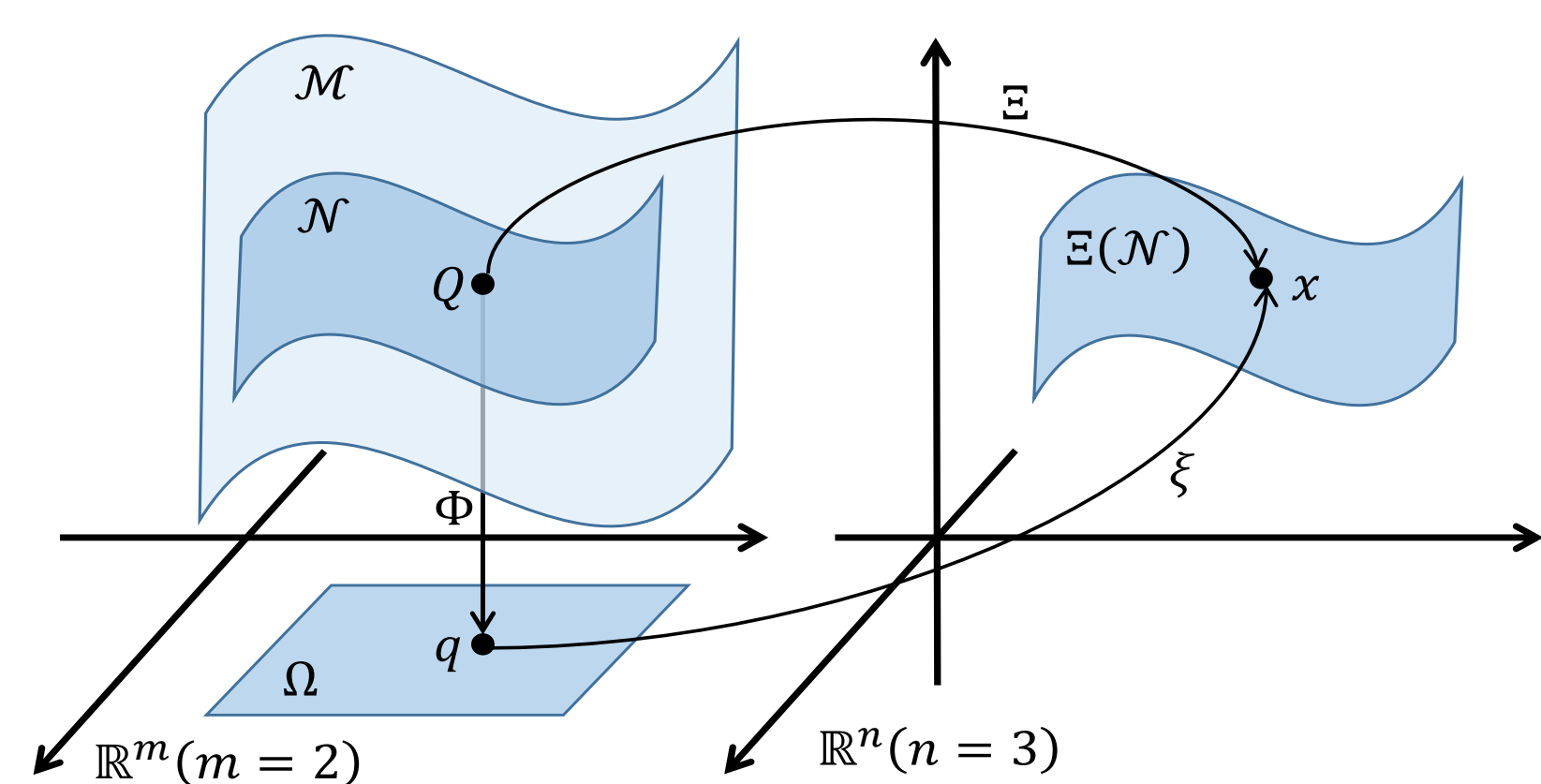
$$dz = f(z)dt + \mathcal{N}(0, 2D(z)dt),$$

with  $f, D$  satisfying certain conditions.

Ability to handle SG noise  $\tilde{f}(z) = f(z) + \mathcal{N}(0, B(z))$ :

$$dz = \tilde{f}(z)dt + \mathcal{N}(0, 2D(z)dt - B(z)dt^2).$$

### Riemann Manifold $\mathcal{M}$



$(\mathcal{N}, \Phi)$ : local coordinate system.

$\Xi$ : (isometric) embedding.  $\xi \triangleq \Xi \circ \Phi^{-1}$ .

$G(q)$ : Riemann metric tensor.

Distribution on  $\mathcal{M}$ :  $\pi_{\mathcal{H}}(x)|_{x=\xi(q)} = \pi(q)/\sqrt{|G(q)|}$ ,

$\pi(q)$ : in the coordinate space;

$\pi_{\mathcal{H}}(x)$ : in the embedded space.

## DYNAMICS CONSTRUCTION

Design novel dynamics in coordinate space by the complete recipe, so that:

- the stationary distribution is desired;
- suitable for 2nd-order integrators.

### SGGMC

Augment with momentum  $p \in \mathbb{R}^m$ :  $z = (q, p)$ .

$$\begin{cases} dq = G^{-1}p dt \\ dp = -\nabla U dt - (1/2)\nabla \log |G| dt - M^T C M G^{-1} p dt \\ \quad - (1/2)\nabla [p^T G^{-1} p] dt + \mathcal{N}(0, 2M^T C M dt) \end{cases}$$

$M(q)_{ij} \triangleq \partial \xi_i(q) / \partial q_j$ ; choose  $C_{n \times n}$  pos. def.

### gSGNHT

$z = (q, p, \xi)$ . Thermostats  $\xi \in \mathbb{R}$ : adaptive  $C$ .

## 2ND-ORDER INTEGRATORS

### Simulate in the Embedded Space

- to release global coordinates requirement;
- $(q, p) \rightarrow (x, v)$ . ( $v$  also momentum)

### Symmetric Splitting Integrator

- Guaranteed to be 2nd-order (Chen et al., 2015).
- Split the dynamics into parts and solve each in closed form:

$$A: dq = G^{-1}p dt, dp = -(1/2)\nabla [p^T G^{-1}p] dt$$

$$B: dp = -M^T C M G^{-1} p dt,$$

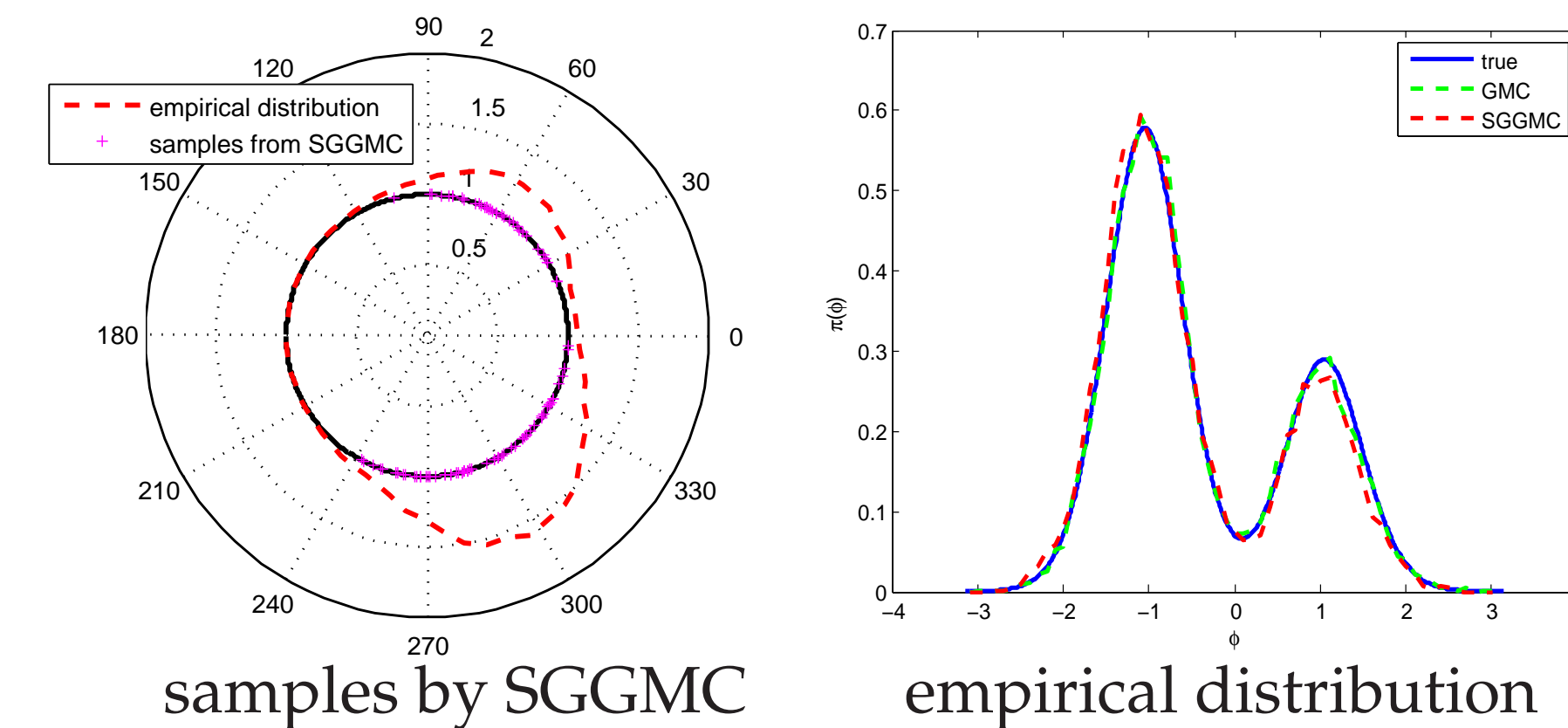
$$O: dp = -\nabla U(q) dt - (1/2)\nabla \log |G| dt + \mathcal{N}(0, 2M^T C M dt).$$

- Simulate the whole dynamics: "ABOBA".

A, B:  $\varepsilon/2$ ; O:  $\varepsilon$ ; use the closed-form solutions.

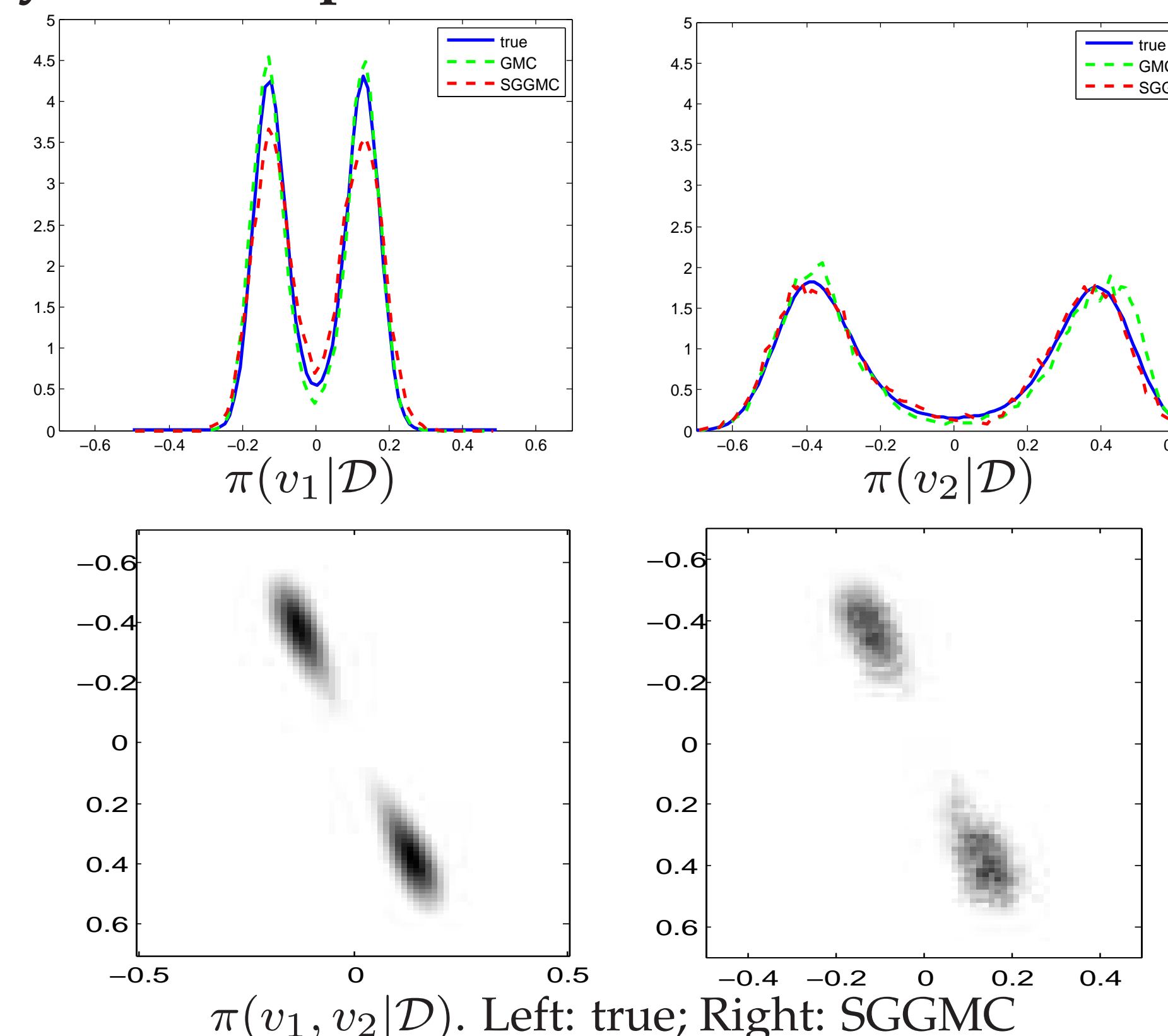
## SYNTHETIC EXPERIMENTS

### Toy Experiment



Sample  $x \in \mathbb{S}^1$ ,  $U(x) = -\log(e^{5\mu_1^T x} + 2e^{5\mu_2^T x})$ . Known gradient noise: artificially added.

### Synthetic Experiment

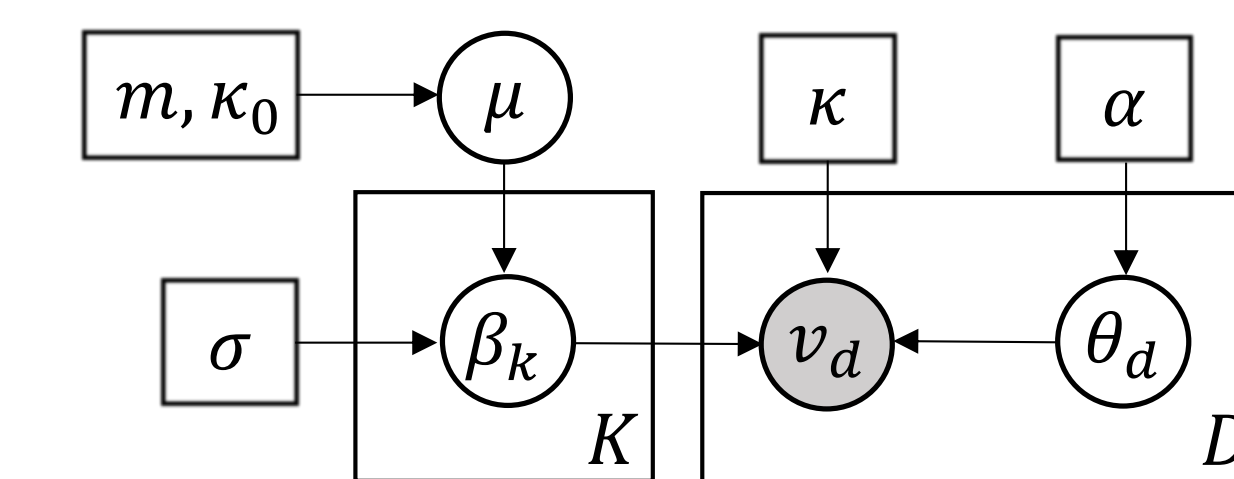


Inference for mixture of vMF:

$$\begin{aligned} \pi(v_1) &= \text{vMF}(v_1|e_1, \kappa_1), \quad \pi(v_2) = \text{vMF}(v_2|e_1, \kappa_2), \\ \pi(x_i|v_1, v_2) &\propto \text{vMF}(x_i|v_1, \kappa_x) + \text{vMF}(x_i|v_2, \kappa_x), \\ \text{with } \mu &\triangleq \frac{v_1 + v_2}{\|v_1 + v_2\|}. \end{aligned}$$

Synthetic data: drawn by GMC.

## REAL-WORLD EXPERIMENTS



Spherical Admixture Model (Reisinger et al., 2014): a topic model for spherical data  $v_d \in \mathbb{S}^{V-1}$ .  $\beta, \mu \in \mathbb{S}^{V-1}$ ,  $\theta \in \Delta^{K-1}$ .

- Inference by SGGMC/gSGNHT: directly sample from the posterior  $\pi(\beta|v)$ .

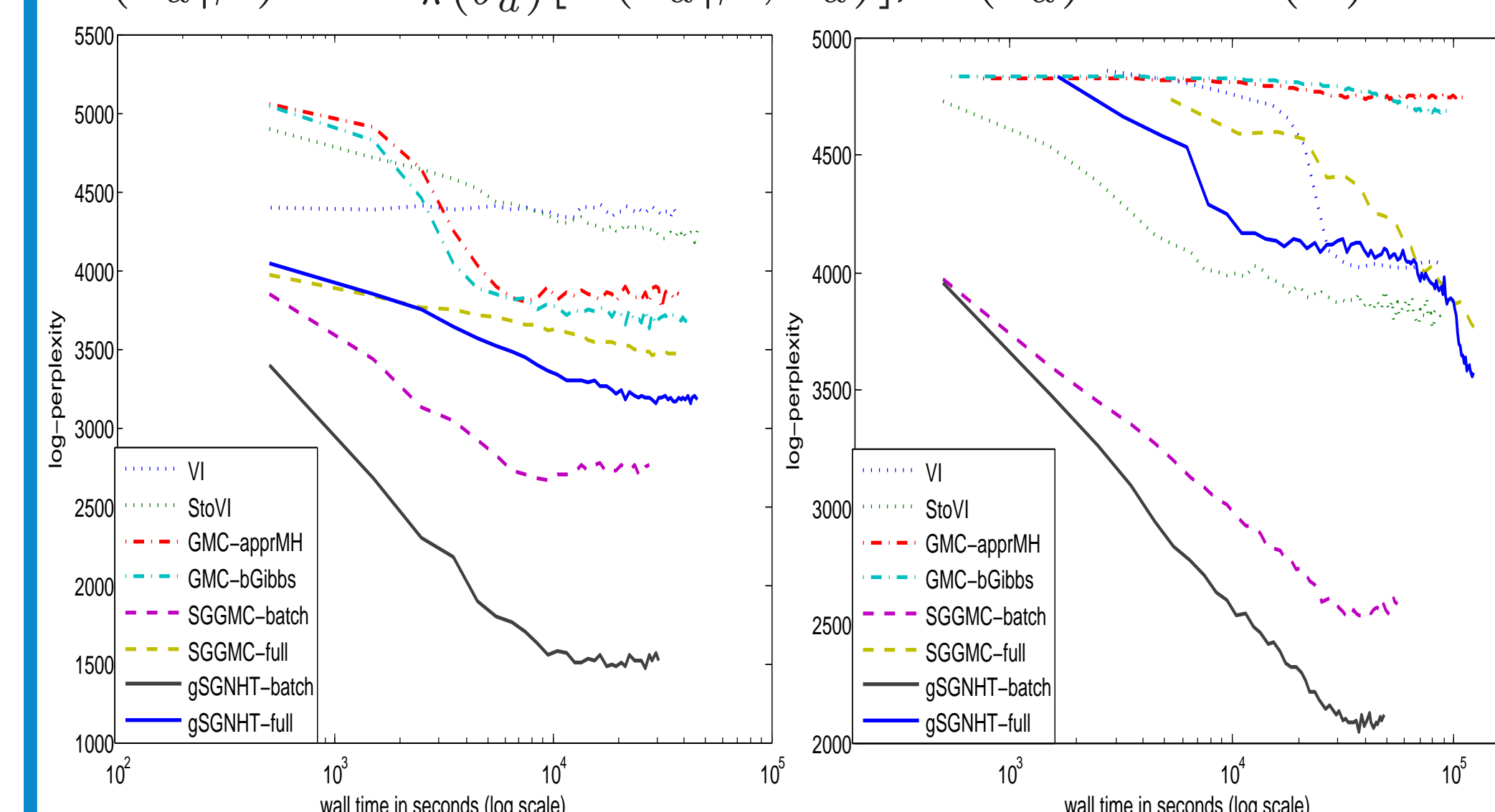
- $\mu$  can be collapsed analytically.
- Monte-Carlo and mini-batch estimated (doubly stochastic) gradient:  $-\nabla_{\beta} \log \pi(\beta|v)$   
 $= -\mathbb{E}_{\pi(\theta|\beta, v)} \left[ \underbrace{\nabla_{\beta} \log \pi(v, \beta, \theta)}_{\text{closed-form known;}} \right]$   
drawn by GMC estimated by mini-batch

- Evaluation:

log-perplexity (the lower the better) along time:

$$-(1/|\mathcal{D}'|) \sum_{d \in \mathcal{D}'} \log \left[ (1/M) \sum_{m=1}^M \pi(v_d|\beta^{(m)}) \right],$$

$$\pi(v_d|\beta) = \mathbb{E}_{\pi(\theta_d)} [\pi(v_d|\beta, \theta_d)], \quad \pi(\theta_d) = \text{Dir}(\alpha).$$



small dataset (1,666)

large dataset (150,000)

- Observations

- SGGMC/gSGNHT: most accurate and fast; more salient on the larger dataset.
- GMC-apprMH/GMC-bGibbs: more accurate than VI/StoVI on the small dataset but too slow on the large one.
- VI/StoVI are blocked.



← Main paper

Appendix, codes, data →

