



# Stein Variational Gradient Descent with Matrix-Valued Kernels

Dilin Wang\*, Ziyang Tang\*, Chandrajit Bajaj, Qiang Liu (\*Equal contribution)  
Department of Computer Science · The University of Texas at Austin

## Overview

### Stein Variational Gradient Descent (SVGD) (Liu and Wang, 2016)

- Given a target distribution  $p$ , find  $q := \{x_i\}_{i=1}^n$  to approximate  $p$ .
- Iteratively move  $q$  towards the target  $p$  by minimizing KL:

$$\min_{\phi \in \mathcal{F}} \text{KL}(q_{[\epsilon\phi]} || p) - \text{KL}(q || p), \text{ with,}$$

$q_{[\epsilon\phi]}$  is the empirical distribution of  $\{x'_i = x_i + \underbrace{\epsilon\phi(x_i)}_{\text{incremental updates}}\}_{i=1}^n$ .

- Closed-form updates by taking  $\mathcal{F}$  to be the unit ball of a RKHS.

### Key limitations:

- Scalar-valued kernels; cannot encode potential correlations between different coordinates.

### Our work:

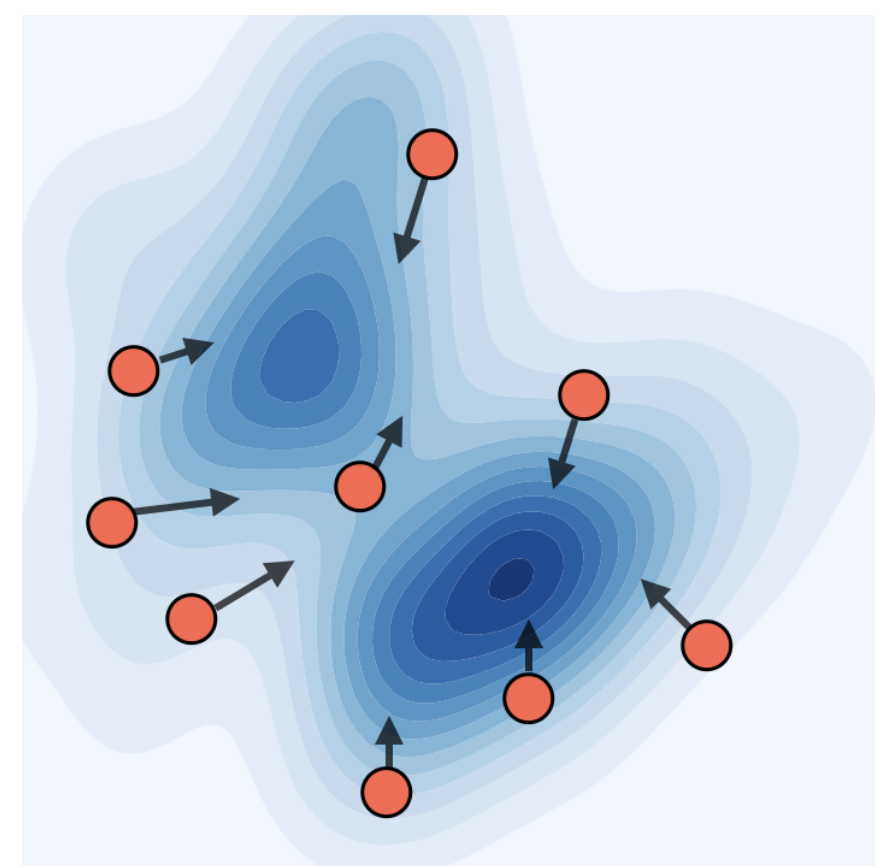
- A generalization of SVGD with matrix-valued kernels.
- Improve SVGD by leveraging geometric information.

## SVGD with scalar-valued kernels

- Finding the optimal transformation  $\phi^*$  such that  $\text{KL}(q||p)$  decreases as fast as possible:

$$\phi^* = \arg \max_{\phi \in \mathcal{F}} \left\{ -\frac{d}{d\epsilon} \text{KL}(q_{[\epsilon\phi]} || p) \Big|_{\epsilon=0} \right\},$$

$$= \arg \max_{\phi \in \mathcal{F}} \left\{ \mathbb{E}_{x \sim q} [\mathcal{P}^\top \phi(x)] \right\},$$



where  $\underbrace{\mathcal{P}^\top \phi(x) := \nabla_x \log p(x)^\top \phi(x) + \nabla_x^\top \phi(x)}_{\text{Stein operator}}$ .

- Taking  $\mathcal{F}$  to be the unit ball of a RKHS  $\mathcal{H}_k^d := \overbrace{\mathcal{H}_k \times \dots \times \mathcal{H}_k}^d$  associated with kernel  $k(x, x')$  (e.g.  $k(x, x') = \exp(-\frac{(x-x')^2}{2h^2})$ ),  
 $\phi_k^*(\cdot) \propto \mathbb{E}_{x \sim q} [k(\cdot, x) \nabla_x \log p(x) + \nabla_x k(\cdot, x)]$ .

## Why matrix-valued kernels

- Improve SVGD with preconditioning

$$x_i = x_i + \epsilon \underbrace{P}_{\text{preconditioning matrix}} \phi_k^*(x_i)$$

A unified approach?

- Matrix-valued kernels

$$\underbrace{K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^{d \times d}}_{\text{Matrix-valued}} \quad \underbrace{k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^1}_{\text{scalar-valued}}.$$

### Example: Separable kernels

·  $K(x, x') = k(x, x')B$ , where  $B \in \mathbb{R}^{d \times d}$ . It is easy to see that the corresponding optimal SVGD update is

$$\phi_K^*(\cdot) = B\phi_k^*(\cdot).$$

## SVGD with Matrix-valued Kernels

- Let  $\mathcal{F}$  be a unit ball of  $\mathcal{H}_K$  with matrix-valued kernels, we have

$$\phi_K^*(\cdot) \propto \mathbb{E}_{x \sim q} [K(\cdot, x) \nabla_x \log p(x) + K(\cdot, x) \nabla_x],$$

$$\text{with, } \underbrace{(K(\cdot, x) \nabla_x)_\ell}_{\ell\text{-th element}} = \sum_{m=1}^d \nabla_{x^m} K_{\ell, m}(\cdot, x).$$

### Main Algorithm

- Initialize** a set of particles  $\{x_i^0\}_{i=1}^n$ ;

- While** NOT converge; do

$$x_i \leftarrow x_i + \frac{\epsilon}{n} \sum_{j=1}^n [K(x_i, x_j) \nabla_{x_j} \log p(x_j) + K(x_i, x_j) \nabla_{x_j}].$$

## Connections with other variants of SVGD

- Vanilla SVGD Liu and Wang, 16

$$K(x, x') = k(x, x')I.$$

- Graphical SVGD Wang et al., 18, Zhuo et al., 18

$$K(x, x') = \text{diag} \left( \left\{ k_\ell(x, x') \right\}_{\ell=1}^d \right).$$

- Gradient-free SVGD Han et al., 18:

$$K(x, x') = k(x, x')w(x)w(x')I, \quad w(x) = \frac{\rho(x)}{p(x)}.$$

## Practical Choices of Matrix-Valued Kernels

### Matrix-SVGD (average)

- Constant preconditioning matrix with RBF kernel:

$$K_Q(x, x') = Q^{-1} \exp \left( -\frac{1}{2h} \|x - x'\|_Q^2 \right),$$

with  $\|x - x'\|_Q^2 = (x - x')^\top Q (x - x')$ .

- $Q = \text{Avg}[H(x_1), \dots, H(x_n)]$ , e.g.  $H$ : Hessian.

### Matrix SVGD (mixture)

- Mixture preconditioning matrix with RBF kernel:

$$K(x, x') = \sum_{\ell=1}^m K_{Q_\ell}(x, x')w_\ell(x)w_\ell(x'),$$

where  $Q_\ell$  as preconditioning of anchor point  $z_\ell$ .

- We set the anchor points to be particles themselves.

- We take  $w_\ell(x)$  as the Gaussian mixture probability:

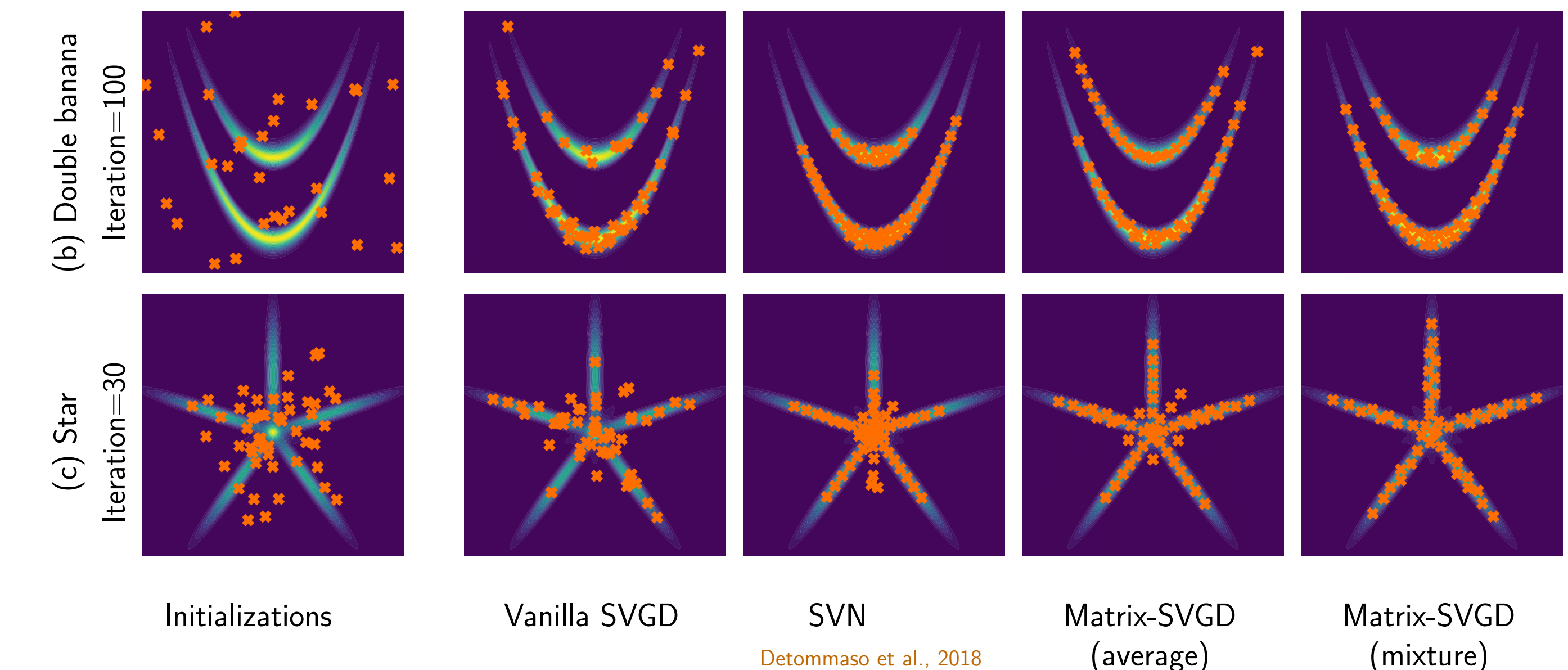
$$w_\ell(x) = \frac{\mathcal{N}(x; z_\ell, Q_\ell^{-1})}{\sum_{\ell=1}^m \mathcal{N}(x; z_\ell, Q_\ell^{-1})},$$

$$\mathcal{N}(x; z_\ell, Q_\ell^{-1}) := \frac{1}{Z_\ell} \exp \left( -\frac{1}{2} \|x - z_\ell\|_{Q_\ell}^2 \right).$$

## Experiments

### Two-Dimensional Toy Examples

- Search the best learning rate for all algorithms exhaustively.



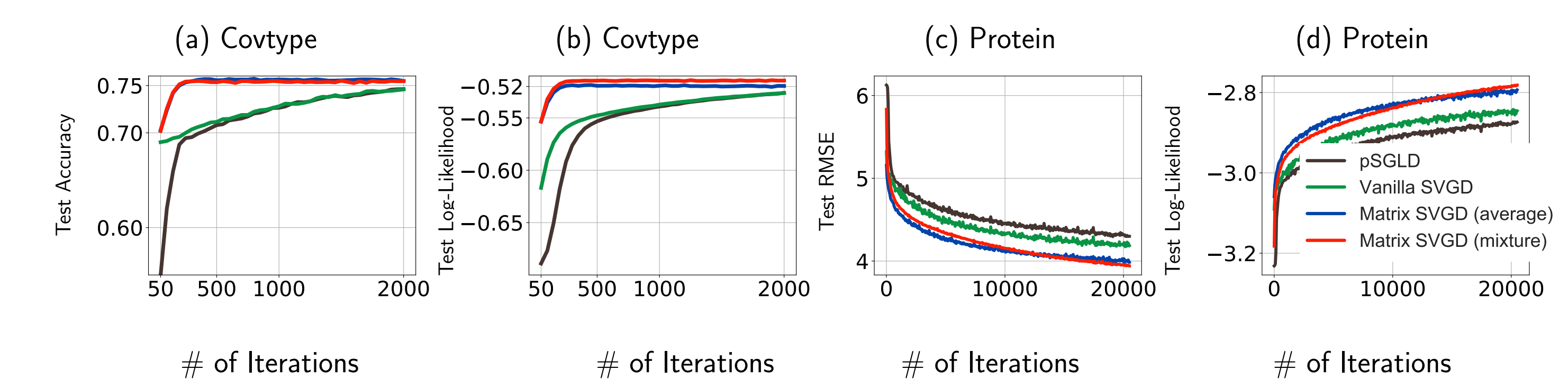
### Bayesian Neural Network Regression

- Matrix-SVGD (mixture) generally performs the best among all algorithms.

Dataset	Vanilla SVGD	Test RMSE		Test Log-Likelihood		
		Vanilla SVGD	Matrix-SVGD (average)	Vanilla SVGD	Matrix-SVGD (average)	Matrix-SVGD (mixture)
Boston	2.785±0.169	2.898±0.184	<b>2.717±0.166</b>	-2.706±0.158	<b>-2.669±0.141</b>	-2.861±0.207
Concrete	5.027±0.116	4.869±0.124	<b>4.721±0.111</b>	<b>-3.064±0.034</b>	-3.150±0.054	-3.207±0.071
Energy	0.889±0.024	<b>0.795±0.025</b>	0.868±0.025	-1.315±0.020	<b>-1.135±0.026</b>	-1.249±0.036
Kin8nm	0.093±0.001	0.092±0.001	<b>0.090±0.001</b>	0.964±0.012	0.956±0.011	<b>0.975±0.011</b>
Naval	0.004±0.000	0.001±0.000	<b>0.000±0.000</b>	4.312±0.087	5.383±0.081	<b>5.639±0.048</b>
Combined	4.088±0.033	4.056±0.033	<b>4.029±0.033</b>	-2.832±0.009	-2.824±0.009	<b>-2.817±0.009</b>
Wine	0.645±0.009	<b>0.637±0.008</b>	<b>0.637±0.009</b>	-0.997±0.019	<b>-0.980±0.016</b>	-0.988±0.018
Protein	4.186±0.017	3.997±0.018	<b>3.852±0.014</b>	-2.846±0.003	-2.796±0.004	<b>-2.755±0.003</b>
Year	8.686±0.010	8.637±0.005	<b>8.594±0.009</b>	-3.577±0.002	-3.569±0.001	<b>-3.561±0.002</b>

### SVGD Acceleration

- (a)-(b) Results of Bayesian Logistic regression on the Covtype dataset.
- (c)-(d) Average test RMSE and log-likelihood vs. training batches on the Protein dataset for Bayesian neural regression.



### Classification With Recurrent Neural Networks

Method	MR	CR	SUBJ	MPQA
	Pang & Lee, 05	Hu & Liu, 04	Pang & Lee, 04	Wiebe et al., 05
SGLD	20.52	18.65	7.66	11.24
pSGLD Li et al., 16	19.75	17.50	6.99	10.80
Vanilla SVGD	19.73	18.07	6.67	<b>10.58</b>
Matrix-SVGD (average)	19.22	17.29	6.76	10.79
Matrix-SVGD (mixture)	<b>19.09</b>	<b>17.13</b>	<b>6.59</b>	10.71

Table: Sentence classification errors measured with four benchmarks.

**Acknowledgment** This work is supported in part by NSF CRII 1830161 and NSF CAREER 1846421. We would like to acknowledge AWS and Google Cloud for their support.

