



# Structured Stochastic Gradient MCMC

Antonios Alexos<sup>1\*</sup> Alex Boyd<sup>2\*</sup> Stephan Mandt<sup>1,2</sup>

<sup>1</sup>Department of Computer Science, <sup>2</sup>Department of Statistics, University of California, Irvine

\* Denotes equal contribution {aalexos, alexjb}@uci.edu

## TL;DR

We propose a new non-parametric variational Langevin-type approximation that makes no parametric assumptions on the posterior distribution. It allows practitioners to impose arbitrary independence structures between parameters, resulting in faster training.

## Preliminaries

- Stochastic gradient Markov chain Monte Carlo (SGMCMC) methods iteratively sample from the posterior distribution:

$$p(\theta|\mathcal{D}) \propto \exp\{-U(\theta)\} \text{ where } U(\theta) = -\log p(\theta, \mathcal{D}).$$

- A classic example is Stochastic Gradient Langevin Dynamics:

$$\theta^{(t+1)} = \theta^{(t)} - \frac{\epsilon_t}{2} \nabla_{\theta} \hat{U}(\theta^{(t)}; \tilde{\mathcal{D}}^{(t)}) + \xi_t, \quad \xi_t \sim \mathcal{N}(0, \epsilon_t I),$$

$$\hat{U}(\theta; \tilde{\mathcal{D}}) = -\frac{|\mathcal{D}|}{|\tilde{\mathcal{D}}|} \log p(\tilde{\mathcal{D}}|\theta) - \log p(\theta), \quad \mathbb{E}_{\tilde{\mathcal{D}}}[\hat{U}(\theta; \tilde{\mathcal{D}})] = U(\theta).$$

- Structured VI best approximates  $p(\theta|\mathcal{D})$  with a restricted distribution  $q(\theta) = \prod_{i=1}^M q_i(\theta_i)$  that assumes  $\theta_i \perp \theta_j$  for  $i \neq j$ .
- This is casted as an optimization problem, minimizing the KL-divergence between them, yielding the following optimal results:

$$\begin{aligned} q^*(\theta) &= \min_q D_{KL}(q(\theta)||p(\theta|\mathcal{D})) = \min_q \mathbb{E}_q \left[ \log \frac{q(\theta)}{p(\theta|\mathcal{D})} \right] \\ &= \exp \left\{ \sum_{i=1}^M \mathbb{E}_{\tilde{\theta}_{-i} \sim q_{-i}} [\log(\theta_i, \tilde{\theta}_{-i}, \mathcal{D})] \right\} \end{aligned}$$

- Typically, parametric assumptions are placed on each  $q_i$  and Coordinate Ascent is performed on the summands of  $q^*(\theta)$ .

## Structured SGMCMC

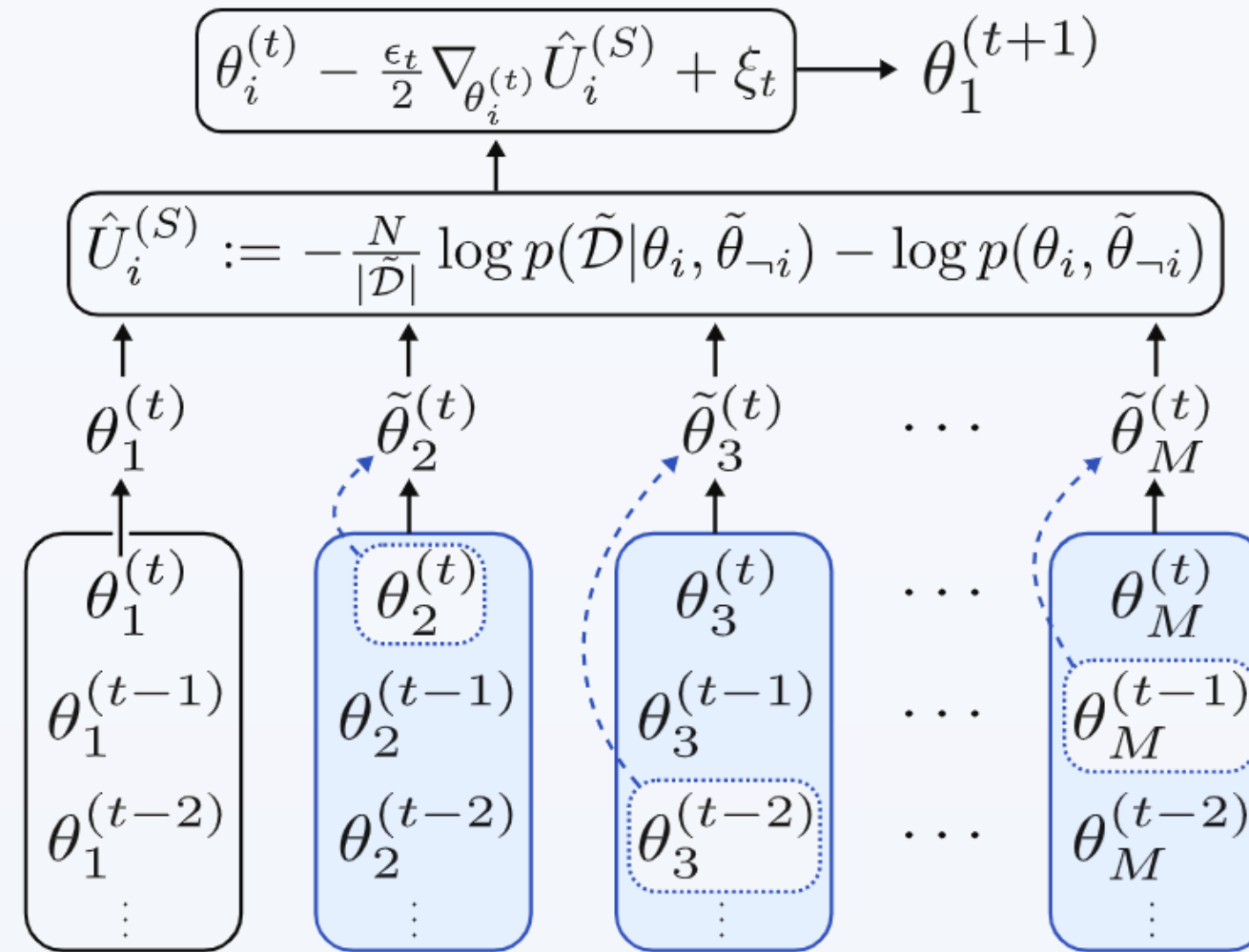
- Instead, we want to sample directly from  $q^*(\theta)$  via SGMCMC and avoid any parametric assumptions on  $q(\theta)$ .
- To do so, we perform SGMCMC with  $\hat{U}(\theta; \tilde{\mathcal{D}})$  replaced by:

$$\begin{aligned} \hat{U}^{(S)}(\theta; \tilde{\mathcal{D}}) &= \sum_{i=1}^M \mathbb{E}_{\tilde{\theta}_{-i} \sim q_{-i}} \hat{U}(\{\theta_i, \tilde{\theta}_{-i}\}; \tilde{\mathcal{D}}), \\ &= \sum_{i=1}^M \mathbb{E}_{\tilde{\theta}_{-i} \sim q_{-i}} \left[ -\frac{|\mathcal{D}|}{|\tilde{\mathcal{D}}|} \log p(\tilde{\mathcal{D}}|\theta_i, \tilde{\theta}_{-i}) - \log p(\theta_i, \tilde{\theta}_{-i}) \right]. \end{aligned}$$

- When generating  $\theta^{(t+1)}$ , we approximate  $\mathbb{E}_{\tilde{\theta}_{-i} \sim q_{-i}}$  with MC samples:

$$\tilde{\theta}_{-i} \sim \prod_{j \neq i} \{\theta_j^{(1)}, \theta_j^{(2)}, \dots, \theta_j^{(t)}\}$$

## Structured SGMCMC Visualized



## Structured Dropout SGMCMC

- Generating a sample for  $\theta$  with Structured SGMCMC requires evaluating  $\hat{U}^{(S)}(\theta; \tilde{\mathcal{D}})$  which requires  $\mathcal{O}(M)$  model forward passes.
- To avoid having computation scaled by the number of parameter groups, we develop a further approximation to sampling from  $q^*(\theta)$ .
- First, we recognize that:

$$\hat{U}^{(S)}(\theta; \tilde{\mathcal{D}}) \equiv M \mathbb{E}_{r \sim \text{Cat}(M^{-1}, \dots, M^{-1})} \mathbb{E}_{\tilde{\theta} \sim q} \hat{U}(r\theta + (1-r)\tilde{\theta}; \tilde{\mathcal{D}})$$

- Using MC samples for the outer expectation breaks scaling issues but leads to sparse gradients ( $\Rightarrow$  not every parameter group being sampled).
- Proposed method replaces  $r \sim \text{Cat}(M^{-1}, \dots, M^{-1})$  with  $r \sim p_{\text{mask}}$  where  $r \in [0,1]^M$  and  $\sum_i r_i > 0$ , which yields a new approximate energy function:

$$\hat{U}^{(Sa)}(\theta; \tilde{\mathcal{D}}) \equiv \frac{M}{\mathbb{E}_{r \sim p_{\text{mask}}} \sum_i r_i} \mathbb{E}_{r \sim p_{\text{mask}}} \mathbb{E}_{\tilde{\theta} \sim q} \hat{U}(r\theta + (1-r)\tilde{\theta}; \tilde{\mathcal{D}})$$

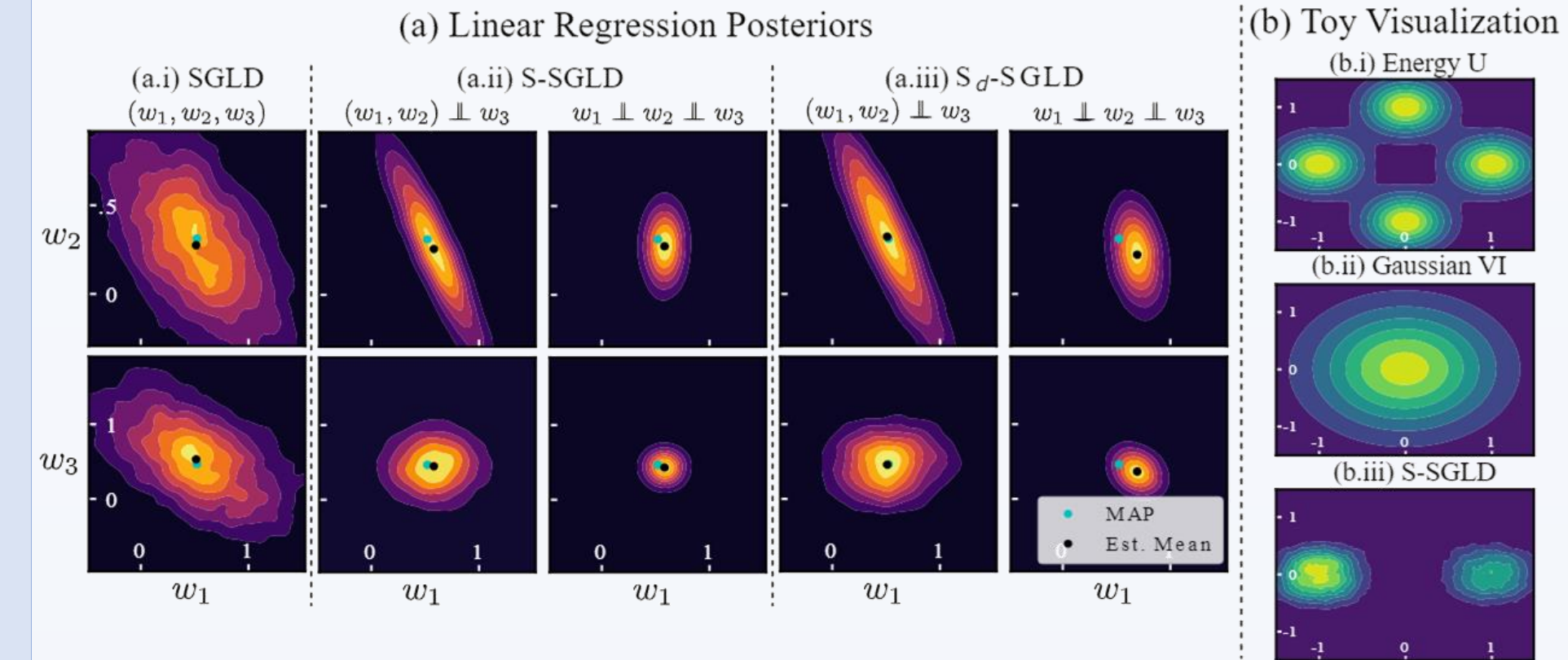
- Example masking distributions:

$$p_{\text{mask}}(r) = \prod_{i=1}^M \text{Unif}(r_i; (0,1)) \quad p_{\text{mask}}(r) = \prod_{i=1}^M \text{Bern}(r_i; \pi) \mathbf{1}(\exists_i r_i = 1)$$

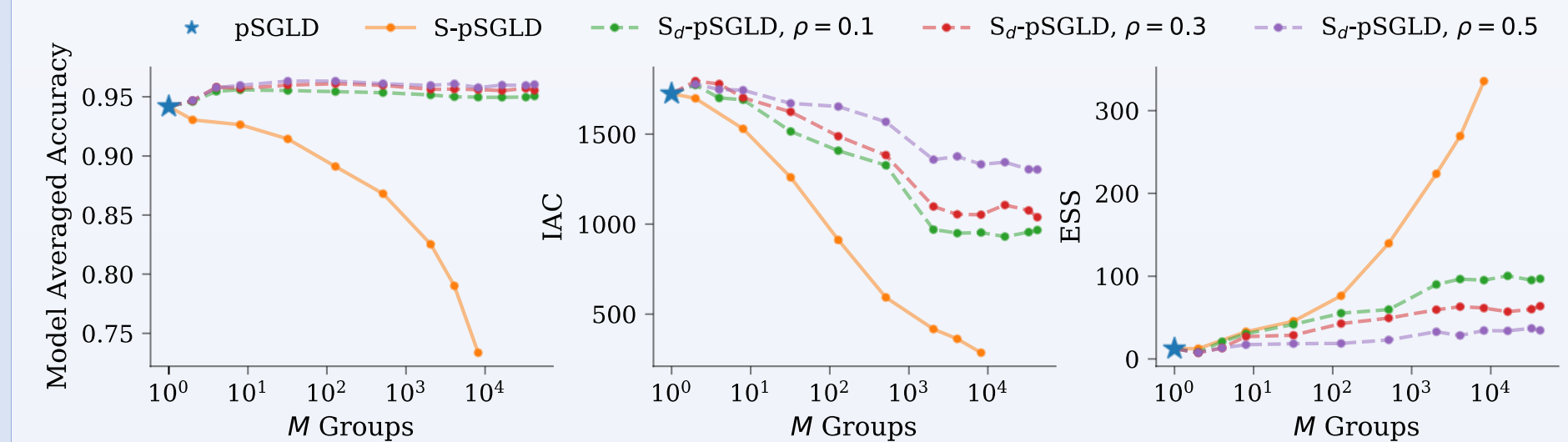
- Structured Dropout SGMCMC can be seen as an interpolation between Structured and Unstructured SGMCMC:

$$r \xrightarrow{a} \text{Cat}(M^{-1}, \dots, M^{-1}) \Rightarrow U^{(Sa)} \rightarrow U^{(S)} \quad r \xrightarrow{a} \{1\}^M \Rightarrow U^{(Sa)} \rightarrow U$$

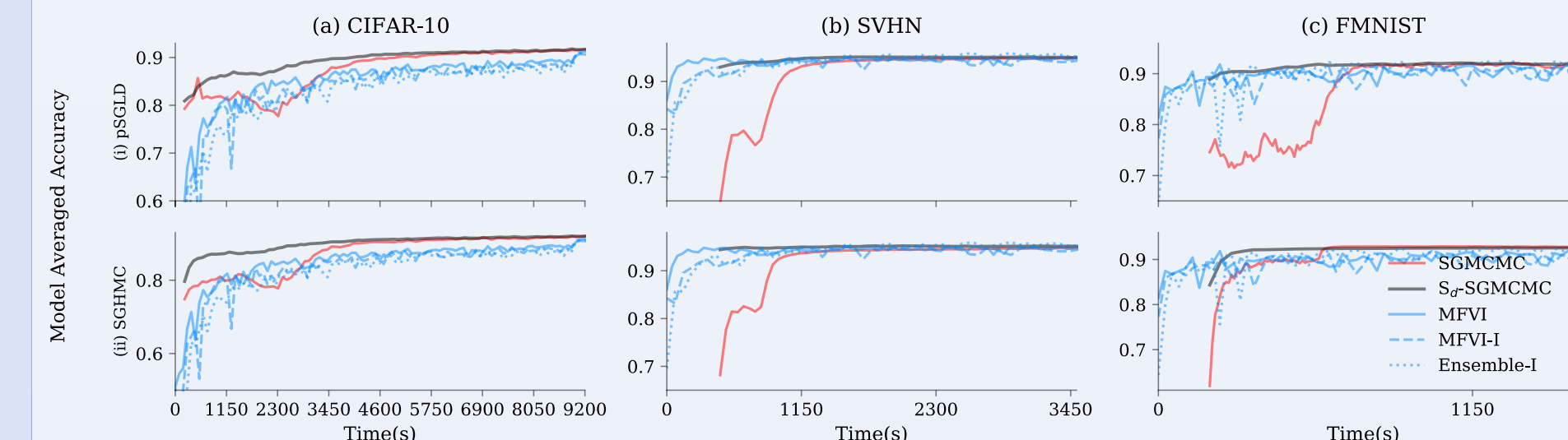
## Experimental Visualizations



## Independence Amount Investigation



- Shown above are the accuracy, integrated autocorrelation time (IAC), and effective sample size (ESS) as a function of the number of independent parameter groups for a NN trained on MNIST.
- More broken correlations equals improvement on IAC, ESS.
- Performance of  $S_r$ -SGMCMC improves compared to S-SGMCMC.



- Accuracy vs. wall-clock time for two base SGMCMC methods,  $S_r$ -SGLD and some VI baselines.  $S_r$ -SGMCMC methods converge much faster and sometimes they outperform the VI baselines.