

Entropy-MCMC: Sampling from Flat Basins with Ease

Bolian Li, Ruqi Zhang

Department of Computer Science
Purdue University



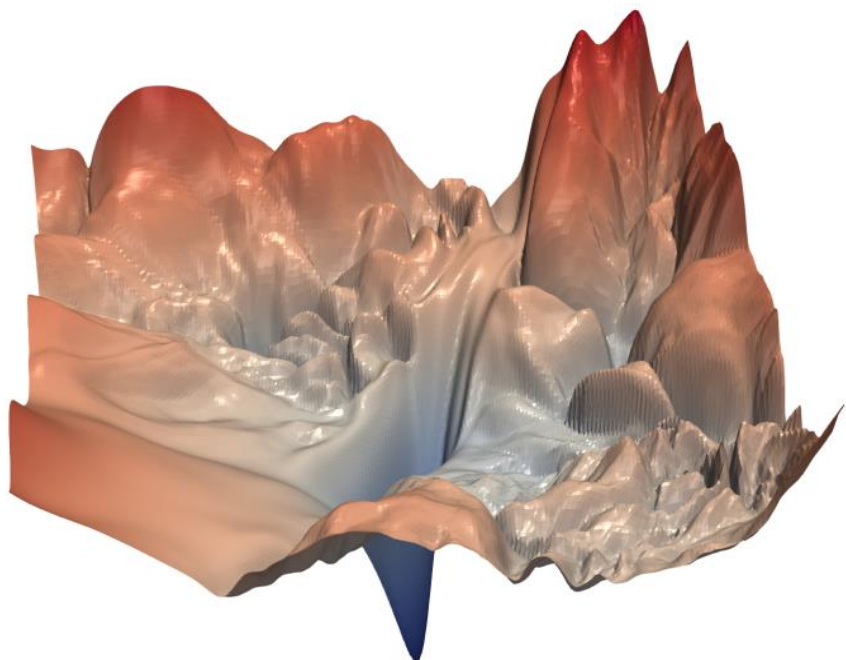
ICLR



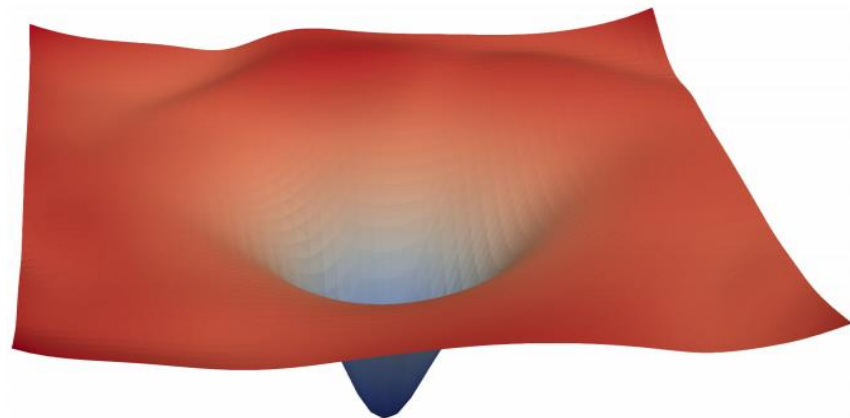
Scan for Paper

Introduction: Loss/Energy Landscape

- Empirical observation: Flat minima generalize better.^[1]



(a) without skip connections

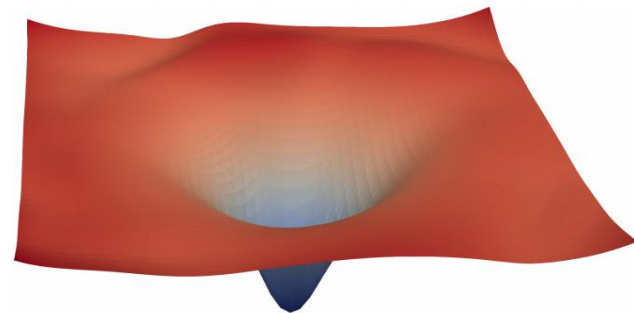


(b) with skip connections

[1] Keskar et al. On large-batch training for deep learning: Generalization gap and sharp minima. ICLR 2017.

[2] Li et al. Visualizing the Loss Landscape of Neural Nets. NeurIPS 2018.

Introduction: Motivation



- Energy landscape of DNNs is highly **multi-modal**.
- Not practical to sample from all modes.
- Flat modes **generalize** better.
- No MCMC methods consider flat minima before.

Preliminaries

Local entropy^[3]:

$$\mathcal{F}(\boldsymbol{\theta}; \eta) = \log \int_{\Theta} \exp \left\{ -f(\boldsymbol{\theta}') - \frac{1}{2\eta} \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|^2 \right\} d\boldsymbol{\theta}'$$

- Averaged energy within a **local region**.
- High local entropy indicates flat regions with low energy values.
- The main objective of Entropy-MCMC.

Stochastic gradient Langevin dynamics (SGLD)^[4]:

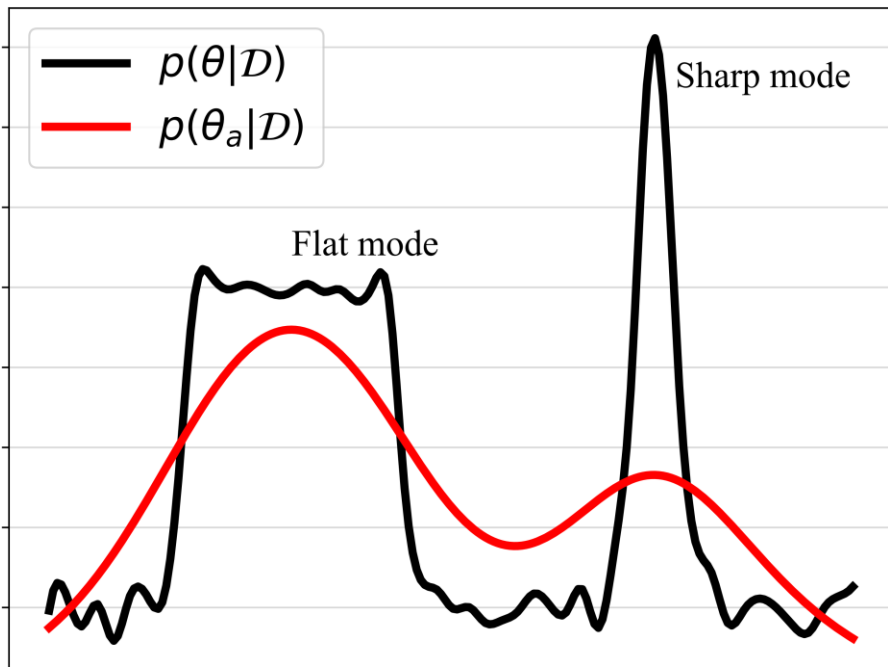
- A standard MCMC algorithm.
- The backbone of Entropy-MCMC implementation.

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} U_{\Xi}(\boldsymbol{\theta}) + \sqrt{2\alpha} \cdot \boldsymbol{\epsilon}$$

[3] Baldassi et al. Subdominant dense clusters allow for simple learning and high computational performance in neural networks with discrete synapses. Physical review letters, 2015.

[4] Welling et al. Bayesian learning via stochastic gradient Langevin dynamics. ICML 2011.

Method: Flat Posterior



- Original posterior: multi-modal, **hard** to sample from

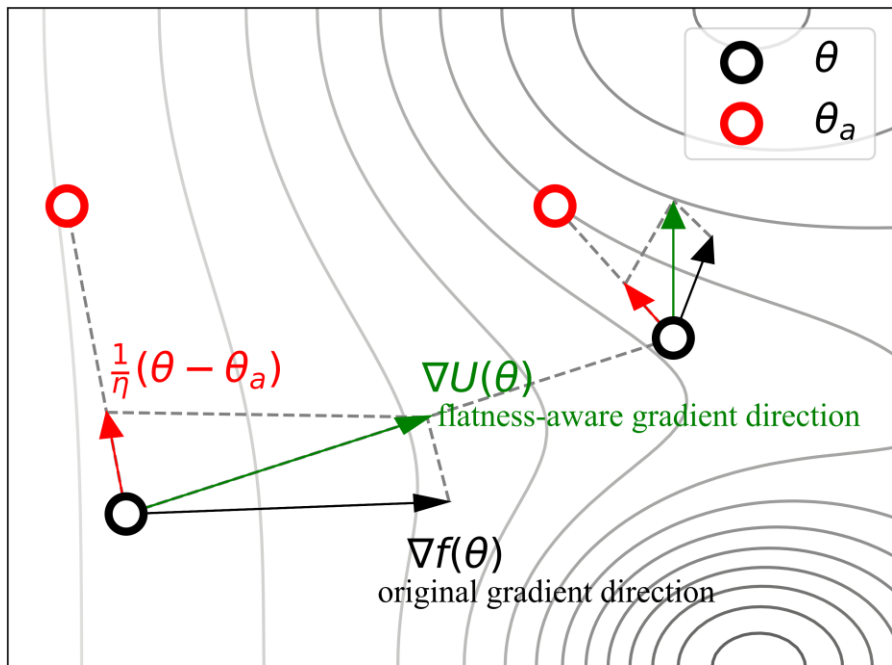
$$p(\theta|\mathcal{D}) \propto \exp(-f(\theta))$$

- Flat posterior: fewer modes, smooth, **easy** to sample from

$$p(\theta_a|\mathcal{D}) \propto \exp \mathcal{F}(\theta_a; \eta) = \int_{\Theta} \exp \left\{ -f(\theta) - \frac{1}{2\eta} \|\theta - \theta_a\|^2 \right\} d\theta$$

- Flat posterior is computed by the **local entropy**.

Method: Sampling



- An auxiliary variable θ_a to eliminate the integral computation

$$p(\theta_a | \mathcal{D}) \propto \exp \mathcal{F}(\theta_a; \eta) = \int_{\Theta} \exp \left\{ -f(\theta) - \frac{1}{2\eta} \|\theta - \theta_a\|^2 \right\} d\theta$$

$$p(\tilde{\theta} | \mathcal{D}) = p(\theta, \theta_a | \mathcal{D}) \propto \exp \left\{ -f(\theta) - \frac{1}{2\eta} \|\theta - \theta_a\|^2 \right\}$$

- For θ , its gradient direction is modified towards flat modes

$$\nabla_{\tilde{\theta}} U(\tilde{\theta}) = \begin{bmatrix} \nabla_{\theta} U(\tilde{\theta}) \\ \nabla_{\theta_a} U(\tilde{\theta}) \end{bmatrix} = \begin{bmatrix} \nabla_{\theta} f(\theta) + \frac{1}{\eta}(\theta - \theta_a) \\ \frac{1}{\eta}(\theta_a - \theta) \end{bmatrix}$$

Method: Sampling

Algorithm 1: Entropy-MCMC

Inputs: The model parameter $\theta \in \Theta$, guiding variable $\theta_a \in \Theta$, and dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$;

Results: Collected samples $\mathcal{S} \subset \Theta$;

$\theta_a \leftarrow \theta, \mathcal{S} \leftarrow \emptyset$;

/ Initialize */*

for *each iteration* **do**

$\Xi \leftarrow$ A mini-batch sampled from \mathcal{D} ;

$U_{\Xi} \leftarrow -\log p(\Xi|\theta) - \log p(\theta) + \frac{1}{2\eta} \|\theta - \theta_a\|^2$;

$\theta \leftarrow \theta - \alpha \nabla_{\theta} U_{\Xi} + \sqrt{2\alpha} \cdot \epsilon_1$;

/ $\epsilon_1, \epsilon_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ */*

$\theta_a \leftarrow \theta_a - \alpha \nabla_{\theta_a} U_{\Xi} + \sqrt{2\alpha} \cdot \epsilon_2$;

if *after burn-in* **then**

$\mathcal{S} \leftarrow \mathcal{S} \cup \{\theta, \theta_a\}$;

/ Collect samples */*

end

end

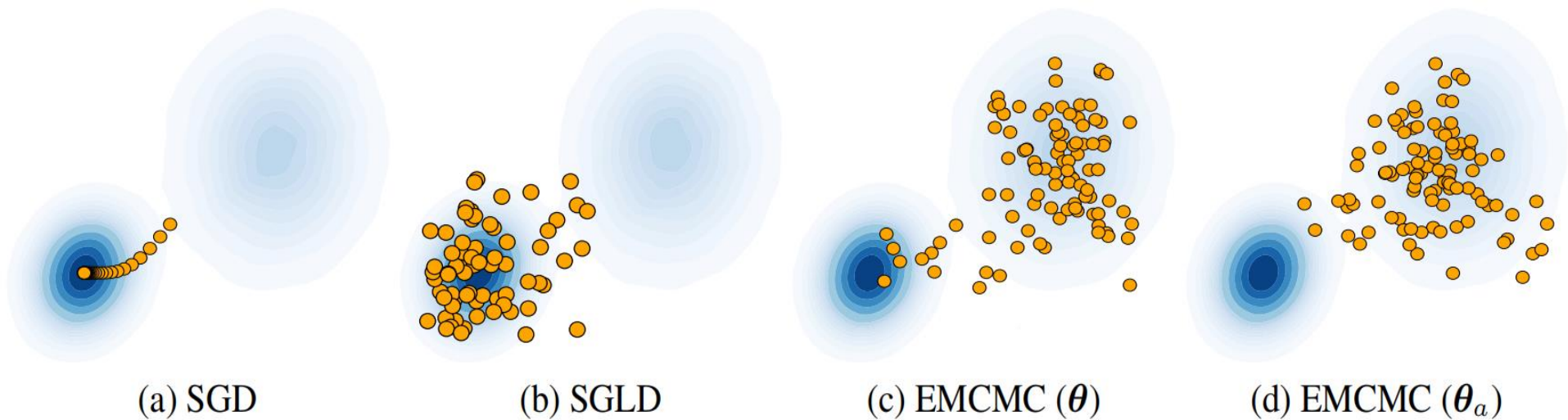
- Inference: Bayesian model averaging

Theoretical Analysis: Convergence Bound

Theorem (informal): Entropy-MCMC converges *faster* than Entropy-SGD and Entropy-SGLD in terms of *2-Wasserstein distance*, due to the removal of *nested Markov chains*.

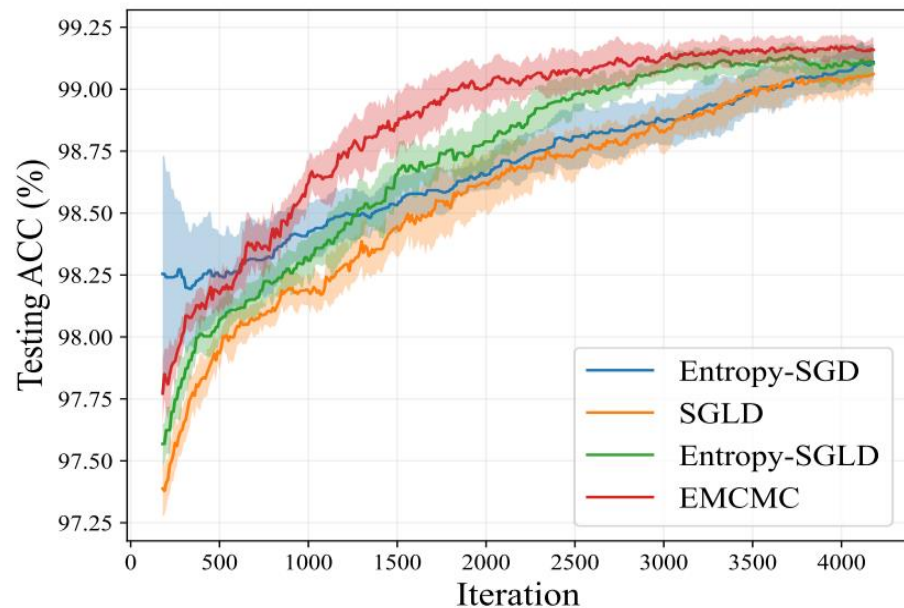
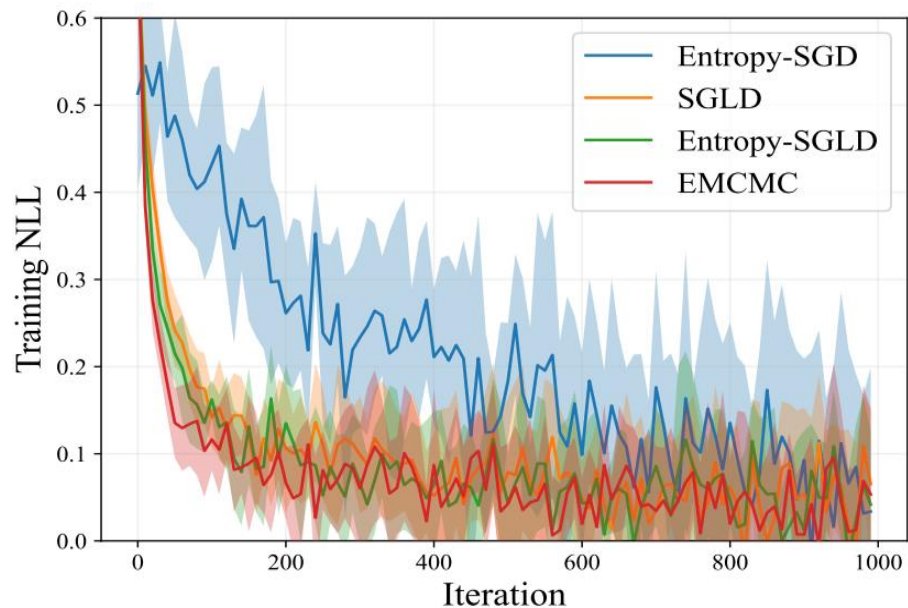
Experiments: Synthetic Examples

- One sharp mode and one flat mode



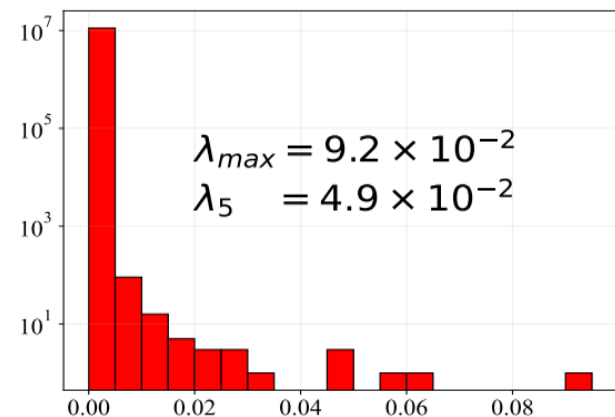
Experiments: Logistic Regression

- Entropy-MCMC converges the fastest

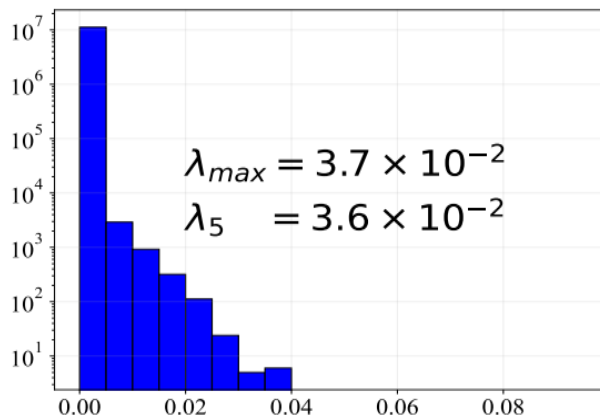


Experiments: Hessian Eigenspectrum

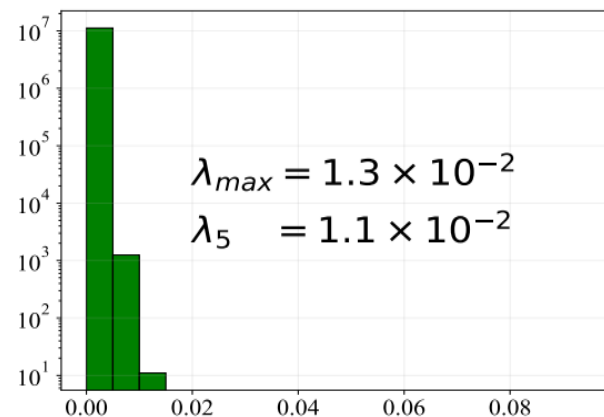
- Lower **Hessian eigenvalues** indicate more flatness



(a) SGD



(b) SGLD



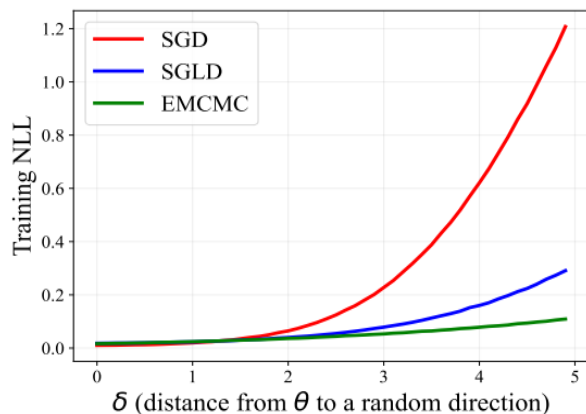
(c) EMCMC

X-axis : eigenvalues

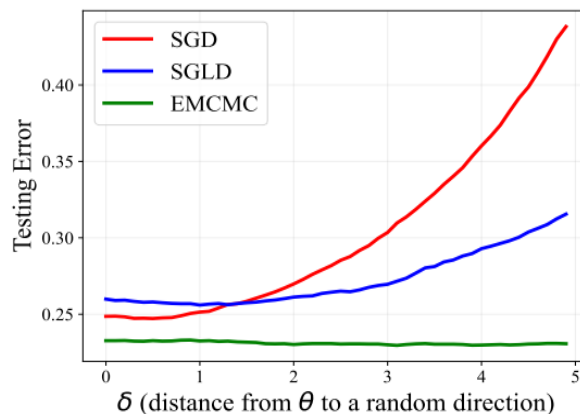
Y-axis : frequency

Experiments: Interpolation

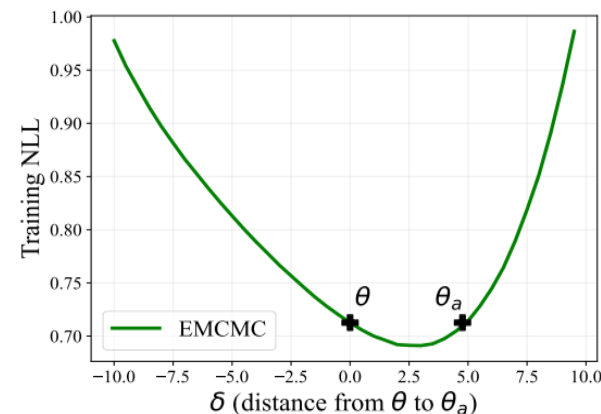
- The mode discovered by Entropy-MCMC is **flatter** than others



(a) $\theta \rightarrow$ Random



(b) $\theta \rightarrow$ Random



(c) $\theta \rightarrow \theta_a$

Experiments: Image Classification

(a) CIFAR10 and CIFAR100

Method	CIFAR10		CIFAR100	
	ACC (%) \uparrow	NLL \downarrow	ACC (%) \uparrow	NLL \downarrow
SGD	94.87 \pm 0.04	0.205 \pm 0.015	76.49 \pm 0.27	0.935 \pm 0.021
Entropy-SGD	95.11 \pm 0.09	0.184 \pm 0.020	77.45 \pm 0.03	0.895 \pm 0.009
SAM	95.25 \pm 0.12	0.166 \pm 0.005	78.41 \pm 0.22	0.876 \pm 0.007
bSAM	95.53 \pm 0.09	0.165 \pm 0.002	78.92 \pm 0.25	0.870 \pm 0.005
SGLD	95.47 \pm 0.11	0.167 \pm 0.011	78.79 \pm 0.35	0.854 \pm 0.031
Entropy-SGLD	94.46 \pm 0.24	0.194 \pm 0.020	77.98 \pm 0.39	0.897 \pm 0.027
EMCMC	95.69 \pm 0.06	0.162 \pm 0.002	79.16 \pm 0.07	0.840 \pm 0.004

(b) Corrupted CIFAR (ACC (%) \uparrow)

Severity	1	2	3	4	5
SGD	88.43	82.43	76.20	67.93	55.81
SGLD	88.61	82.46	76.49	69.19	56.98
EMCMC	88.87	83.27	77.44	70.31	58.17

(c) ImageNet

Metric	NLL \downarrow	Top-1 (%) \uparrow	Top-5 (%) \uparrow
SGD	0.960	76.046	92.776
SGLD	0.921	76.676	93.174
EMCMC	0.895	77.096	93.424

Experiments: OOD detection

- Predictive uncertainty
- Good characterization of posterior leads to good OOD detection

Method	CIFAR10-SVHN		CIFAR100-SVHN	
	AUROC (%) \uparrow	AUPR (%) \uparrow	AUROC (%) \uparrow	AUPR (%) \uparrow
SGD	98.30	99.24	71.96	84.08
Entropy-SGD	98.71	99.37	79.15	86.92
SAM	94.23	95.67	74.56	84.61
SGLD	97.66	98.64	72.51	83.35
Entropy-SGLD	90.07	91.80	71.83	82.89
EMCMC	98.15	99.04	81.14	87.18

Conclusion

1. Sampling from the **flat basins** can improve the generalization of MCMC samples.
2. The proposed **joint posterior** distribution can eliminate the need for integral computation.
3. Entropy-MCMC can effectively find flat modes and achieve promising empirical results.

Thank You!