

# Communication Efficient Stochastic Gradient MCMC for Neural Networks

## Supplementary Material

Chunyuan Li<sup>1</sup>, Changyou Chen<sup>2</sup>, Yunchen Pu<sup>3</sup>, Ricardo Henao<sup>4</sup>, and Lawrence Carin<sup>4</sup>

<sup>1</sup>Microsoft Research, Redmond <sup>2</sup>University at Buffalo, SUNY <sup>3</sup>Facebook <sup>4</sup>Duke University

### A. Algorithm of Downpour pSGLD

---

#### Algorithm 1 Downpour pSGLD

---

```

1: Input:  $\epsilon_t, \beta_1, \lambda, \pi \in \mathbb{N}$ 
2: Output:  $\{\theta_l\}_{l=1}^L$ 
3: Initialize:  $t^{(p)} = 0, l = 0, \nu^{(p)} = 0$ 
    $\tilde{\theta} \sim \mathcal{N}(0, \mathbf{I}), \theta^{(p)} = \tilde{\theta}$ 
4: while maximum time not reached do
5:   % Estimate gradient from  $\mathcal{D}_M$ 
6:    $\tilde{\mathbf{f}}_t^{(p)} = \nabla \tilde{U}(\theta_t^{(p)})$ 
7:   % Preconditioning
8:    $\mathbf{v}_t^{(p)} \leftarrow \beta_1 \mathbf{v}_{t-1} + (1 - \beta_1) \tilde{\mathbf{f}}_t \odot \tilde{\mathbf{f}}_t$ 
9:    $\mathbf{G}_t^{(p)} \leftarrow \text{diag} \left( \mathbf{1} \oslash (\lambda \mathbf{1} + \sqrt{\mathbf{v}_t^{(p)}}) \right)$ 
10:  % Parameter update with pSGLD
11:   $\xi_t^{(p)} \sim \mathcal{N}(0, \mathbf{I})$ 
12:   $\theta_{t+1}^{(p)} \leftarrow \theta_t^{(p)} - \epsilon_t \mathbf{G}_t^{(p)} \tilde{\mathbf{f}}_t^{(p)} + \sqrt{\epsilon_t \mathbf{G}_t^{(p)}} \xi_t^{(p)}$ 
13:   $\nu_{t+1}^{(p)} \leftarrow \nu_t^{(p)} - \epsilon_t \mathbf{G}_t^{(p)} \tilde{\mathbf{f}}_t^{(p)} + \sqrt{\epsilon_t \mathbf{G}_t^{(p)}} \xi_t^{(p)}$ 
14:   $t^{(p)} \leftarrow t^{(p)} + 1$ 
15:  %  $\pi$ -Downpour Communication
16:  if  $t^{(p)}$  divide  $\pi$  then
17:     $\theta_{l+1} \leftarrow \theta_l + \nu_{t+1}^{(p)}$ 
18:     $\theta_{t+1}^{(p)} \leftarrow \theta_{l+1} \quad \nu_{t+1}^{(p)} \leftarrow 0$ 
19:     $l = l + 1$ 
20:  end if
21: end while
```

---

Preconditioned SGLD (pSGLD) (Li *et al.* 2016) specifies the general dynamics of SG-MCMC as  $\mathbf{z} = \theta$ ,  $H(\theta) = U(\theta)$ ,  $W(\theta) = \mathbf{G} = G(\theta)$ ,  $Q(\theta) = \mathbf{0}$ , and  $\Gamma(\theta)$  is a vector with its  $i$ -th element as  $\sum_j \frac{\partial}{\partial \theta_j} (G_{ij})$ . pSGLD utilizes magnitudes of recent gradients to construct a diagonal preconditioner  $\mathbf{G}$ . The downpour pSGLD is shown in Algorithm 1. There are two tuning parameters in pSGLD:  $\lambda$  controls the extremes of the curvature in the preconditioner (default  $\lambda=10^{-5}$ ), and  $\beta_1$  balances the weights of historical and current gradients. We use a default value of  $\beta_1=0.99$ , to construct an exponentially decaying sequence.

### B. MSE bounds for Communication Protocols

We extend the proofs in (Chen *et al.* 2015) to study the impact of the proposed periodic communication protocols on the MSE bounds for accelerating SG-MCMC algorithms. We follow the same assumptions of (Chen *et al.* 2015; Teh *et al.* 2016; Vollmer *et al.* 2015).

The MSE bound for standard SG-MCMC is:

$$\mathbb{E} \left( \hat{\phi}_L - \bar{\phi} \right)^2 \leq \mathcal{B}_0 \triangleq C(\mathcal{E}_1 + \mathcal{E}_2), \text{ with} \quad (1)$$

$$\mathcal{E}_1 = \sum_l \frac{\epsilon_l^2}{S_L^2} \mathbb{E} \|\Delta V_l\|^2 \quad \text{and} \quad \mathcal{E}_2 = \frac{1}{S_L} + \frac{(\sum_l \epsilon_l^2)^2}{S_L^2}, \quad (2)$$

#### B2. Proof for Downpour Protocol

**Proof** To simplify the proof, we note that the master collects a sample every  $\pi$  iterations for each worker, this corresponds to thinning the worker chains, with thinning interval  $\pi$ . According to Corollary 2 of (Li *et al.* 2016), this is equivalent to collecting samples from the chains with step-size  $\epsilon_l^{(p)} = \sum_{t=t^{(p)}-\pi}^{t^{(p)}} \epsilon_t$ . As a result, the algorithm is transformed to a simpler algorithm, in which the chains communicate with the master instantly.

Compared with standard SG-MCMC, there is an additional error associated with the gradient approximation, *i.e.*,  $\theta_{l+1}$  is obtained with a stochastic gradient  $\tilde{\mathbf{f}}_{l-\pi_l}$  evaluated on “old” parameters  $\theta_{l-\pi_l}$  for some integer  $\pi_l$ , instead of  $\theta_l$ . In other words,  $\tilde{\mathbf{f}}_{l-\pi_l}$  has been used to approximate  $\tilde{\mathbf{f}}_l$ , resulting in additional approximation error.

Only the term  $\mathcal{E}_1$  relies on the gradient, in which  $\Delta V_l$  needs to be replaced with

$$\begin{aligned} \Delta \tilde{V}_l &\triangleq \tilde{\mathbf{f}}_{l-\pi_l} - \mathbf{f}_l \\ &= (\tilde{\mathbf{f}}_{l-\pi_l} - \tilde{\mathbf{f}}_l) + (\tilde{\mathbf{f}}_l - \mathbf{f}_l) \\ &= (\tilde{\mathbf{f}}_{l-\pi_l} - \tilde{\mathbf{f}}_l) + \Delta V_l, \end{aligned}$$

where the “old” parameter is at most  $\pi(P-1)$  steps away from the latest parameter, as there are totally  $P$  workers:

$$\pi_l \leq \pi(P-1).$$

As a result, we need to bound the term  $\sum_l \frac{\epsilon_l^2}{S_L^2} \mathbb{E} \|\Delta \tilde{V}_l\|^2$ ,

which gives

$$\begin{aligned} & \sum_l \frac{\epsilon_l^2}{S_L^2} \mathbb{E} \left\| \Delta \tilde{V}_l \right\|^2 \\ & \leq \sum_l \frac{\epsilon_l^2}{S_L^2} \mathbb{E} \left\| \Delta V_l \right\|^2 + \sum_l \frac{\epsilon_l^2}{S_L^2} \mathbb{E} \left\| \tilde{\mathbf{f}}_{l-\pi_l} - \tilde{\mathbf{f}}_l \right\|^2 \\ & \leq \mathcal{E}_1 + D_1 \sum_l \frac{\epsilon_l^2}{S_L^2} \mathbb{E} \left\| \boldsymbol{\theta}_{l-\pi_l} - \boldsymbol{\theta}_l \right\|^2, \end{aligned}$$

where the last inequality follows from the Lipschitz property. Now, using the integrator property as in (Chen *et al.* 2015), we have  $\mathbb{E} \left\| \boldsymbol{\theta}_l - \boldsymbol{\theta}_{l+\delta} \right\| = O(\delta \epsilon_l^K)$  for a  $K$ -th-order integrator. Because we are using a 1st-order integrator ( $K = 1$ ), we have

$$\begin{aligned} \mathbb{E} \left\| \boldsymbol{\theta}_{l-\pi_l} - \boldsymbol{\theta}_l \right\|^2 &= (\mathbb{E} \left\| \boldsymbol{\theta}_{l-\pi_l} - \boldsymbol{\theta}_l \right\|)^2 + \text{Var}(\boldsymbol{\theta}_{l-\pi_l} - \boldsymbol{\theta}_l) \\ &\leq (\mathbb{E} \left\| \boldsymbol{\theta}_{l-\pi_l} - \boldsymbol{\theta}_l \right\|)^2 \leq D_2 \pi_l^2, \end{aligned}$$

for some constant  $D_2$ .

As a result, we have

$$\begin{aligned} \sum_l \frac{\epsilon_l^2}{S_L^2} \mathbb{E} \left\| \Delta \tilde{V}_l \right\|^2 &\leq \mathcal{E}_1 + D_3 \sum_l \frac{\epsilon_l^2}{S_L^2} \pi_l^2 \\ &\leq \mathcal{E}_1 + D_3 \sum_l \frac{\epsilon_l^2}{S_L^2} \pi^2 (P-1)^2. \end{aligned}$$

Substitute the above inequality into the MSE bound, we get the MSE bound for downpour protocol as

$$\begin{aligned} \mathbb{E} \left( \hat{\phi}_L - \bar{\phi} \right)^2 &\leq C_1 (\mathcal{E}_1 + \mathcal{E}_2 + \mathcal{E}_3), \text{ with} \\ \mathcal{E}_3 &= \pi^2 (P-1)^2 \frac{(\sum_l \epsilon_l^2)^2}{S_L^2}, \end{aligned}$$

The proof for bias and variance follow similar derivation. ■

## B2. Proof for Elastic Protocol with $\alpha = 1$

**Proof**

$$\begin{aligned} \mathbb{E} \left( \hat{\phi}_L - \bar{\phi} \right)^2 &= \mathbb{E} \left( \sum_{p=1}^P \frac{S_L^{(p)}}{\sum_p S_L^{(p)}} \left( \hat{\phi}_L^{(p)} - \bar{\phi} \right) \right)^2 \\ &= \sum_{p=1}^P \frac{S_L^{(p)^2}}{(\sum_p S_L^{(p)})^2} \mathbb{E} \left( \hat{\phi}_L^{(p)} - \bar{\phi} \right)^2 \\ &\quad + \sum_{i \neq j} \frac{S_L^{(i)} S_L^{(j)}}{(\sum_p S_L^{(p)})^2} \mathbb{E} \left[ \hat{\phi}_L^{(i)} - \bar{\phi} \right] \mathbb{E} \left[ \hat{\phi}_L^{(j)} - \bar{\phi} \right] \\ &\leq \sum_{p=1}^P \frac{S_L^{(p)^2}}{(\sum_p S_L^{(p)})^2} \mathbb{E} \left( \hat{\phi}_L^{(p)} - \bar{\phi} \right)^2 \\ &\quad + \sum_{i \neq j} \frac{S_L^{(i)} S_L^{(j)}}{(\sum_p S_L^{(p)})^2} \left| \mathbb{E} \hat{\phi}_L^{(i)} - \bar{\phi} \right| \left| \mathbb{E} \hat{\phi}_L^{(j)} - \bar{\phi} \right|. \end{aligned}$$

From (Chen *et al.* 2015), the bias and MSE bounds above are formulated as

$$\begin{aligned} \mathbb{E} \left( \hat{\phi}_L^{(p)} - \bar{\phi} \right)^2 &= O \left( \mathcal{E}_1^{(p)} + \frac{1}{S_L^{(p)}} + \frac{(\sum_l \epsilon_l^{(p)^2})^2}{S_L^{(p)^2}} \right) \\ \left| \mathbb{E} \hat{\phi}_L^{(p)} - \bar{\phi} \right| &= O \left( \frac{1}{S_L^{(p)}} + \frac{\sum_l \epsilon_l^{(p)^2}}{S_L^{(p)}} \right) \end{aligned}$$

Substituting the bounds into the above MSE for elastic protocol, we have

$$\begin{aligned} \mathbb{E} \left( \hat{\phi}_L - \bar{\phi} \right)^2 &\leq \sum_{p=1}^P \frac{S_L^{(p)^2}}{(\sum_p S_L^{(p)})^2} \left( \mathcal{E}_1^{(p)} + \frac{1}{S_L^{(p)}} + \frac{(\sum_l \epsilon_l^{(p)^2})^2}{S_L^{(p)^2}} \right) \\ &\quad + \sum_{i \neq j} \frac{S_L^{(i)} S_L^{(j)}}{(\sum_p S_L^{(p)})^2} \left( \frac{1}{S_L^{(i)}} + \frac{\sum_l \epsilon_l^{(i)^2}}{S_L^{(i)}} \right) \left( \frac{1}{S_L^{(j)}} + \frac{\sum_l \epsilon_l^{(j)^2}}{S_L^{(j)}} \right) \\ &\leq D_4 \left( \sum_{p=1}^P \frac{S_L^{(p)^2}}{(\sum_p S_L^{(p)})^2} \mathcal{E}_1^{(p)} + \frac{1}{\sum_p S_L^{(p)}} + \sum_{p=1}^P \frac{(\sum_l \epsilon_l^{(p)^2})^2}{(\sum_p S_L^{(p)})^2} \right. \\ &\quad \left. + \frac{P^2 - P}{(\sum_p S_L^{(p)})^2} + \sum_{i \neq j} \frac{(\epsilon_l^{(i)^2})(\epsilon_l^{(j)^2})}{(\sum_p S_L^{(p)})^2} \right), \end{aligned}$$

for some constant  $D_4$ .

$$\begin{aligned} &\leq D_5 \left( \sum_{p=1}^P \frac{S_L^{(p)^2}}{(\sum_p S_L^{(p)})^2} \mathcal{E}_1^{(p)} \right. \\ &\quad \left. + \frac{1}{\sum_p S_L^{(p)}} + \sum_{i,j=1}^P \frac{(\sum_l \epsilon_l^{(i)^2} + \sum_l \epsilon_l^{(j)^2})^2}{(\sum_p S_L^{(p)})^2} \right) \end{aligned}$$

The proof for bias and variance follow similar derivation, where the variance (Chen *et al.* 2016) is

$$\begin{aligned} \mathcal{V}_2 &= V_2 \left( \frac{1}{\sum_p S_L^{(p)}} + \sum_{p=1}^P \frac{\sum_l (\epsilon_l^{(p)})^2}{(\sum_p S_L^{(p)})^4} \left( \sum_l \epsilon_l^{(p)} \right)^2 \right) \\ &\leq V_2 \left( \frac{1}{\sum_p S_L^{(p)}} + \sum_{i,j=1}^P \frac{(\sum_l \epsilon_l^{(i)^2} + \sum_l \epsilon_l^{(j)^2})^2}{(\sum_p S_L^{(p)})^2} \right) \end{aligned}$$

■

## C. Additional Experimental Results

**Details on settings** We have focused on the impact of  $\pi$  and  $P$  in the main text. In SG-MCMC algorithms, the hyperparameters for our methods also include mini-batch size, thinning interval, number of burn-in's, step-size and variances of the Gaussian priors. We discuss and specify them in the following.

- **Mini-batch size  $M$ :** Following conventional settings on these datasets, we choose  $M = 100$  throughout the experiments.
- **Burn-in:** To obtain a good initialization for parameter samples from regions of higher probability, "burn-in" is usually used in MCMC to dispose of samples at the beginning of an MCMC run. We choose zero burn-in for MNIST/SVHN, a small burn-in period for WP in our parallel SG-MCMC methods. Because we aim to use multiple workers to explore the space.

- **Thinning interval:** While the proposed periodic protocol can play the role of thinning, we also manually thin the master chain: we collect samples from the master chain every  $R$  seconds in wall-clock time for testing. This is because (i) the same number of model samples is kept in testing with standard SG-MCMC, for fair comparison, and (ii) it can reduce the computational burden during testing. Specifically,  $R = 5$  for MNIST,  $R = 15$  for SVHN, and  $R = 30$  for WP.
- **Step-size:** The default step-size for SGHMC and SGLD is  $1 \times 10^{-2}$ , and  $1 \times 10^{-3}$  for pSGLD.
- **Gaussian priors:** The variance of the Gaussian distribution encapsulates the prior belief of users on how strongly these weights should concentrate around 0. A larger variance in the prior encourages a wider range of weight choices, thus higher uncertainty in learning the weights. We fix  $\sigma^2 = 1$  to focus on the study of communication protocols, and refer readers to (Li *et al.* 2016) for details about the impact of variance in modeling weight uncertainty.

**More Results on FNN** Figure 1 shows the test classification errors for various parallel algorithms.

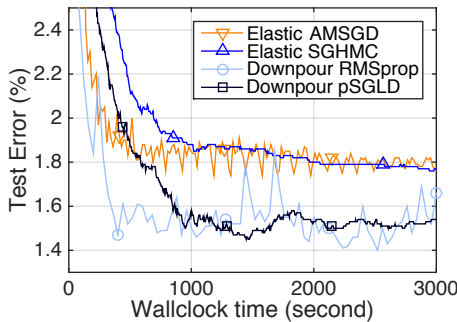


Figure 1: Comparison of Classification errors on FNN.

**More Results on CNN** Figure 2 shows the test classification errors of parallel SG-MCMC algorithms with 1 and 4 workers. Lower classification errors can be achieved when using 4 workers compared to 1 worker, *i.e.*, standard SG-MCMC.

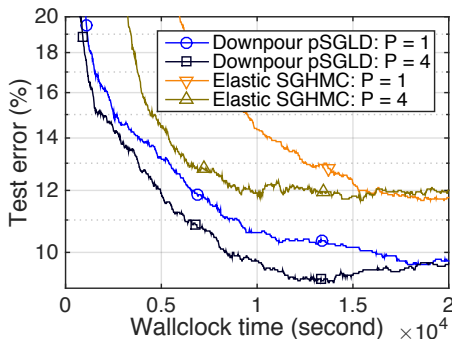


Figure 2: Classification error on CNN.

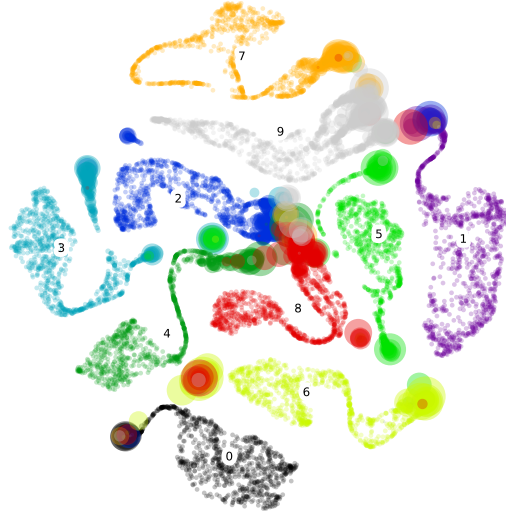


Figure 3:  $t$ -SNE visualization of MNIST.

Table 1: NLL on WP dataset.

Gates		LSTM		GRU	
#Layers		1-layer	2-layer	1-layer	2-layer
Downpour pSGLD	P=1	1.778	1.408	1.705	1.578
	P=4	<b>1.582</b>	<b>1.266</b>	<b>1.548</b>	<b>1.400</b>
Elastic SGHMC	P=1	1.513	1.429	1.608	1.351
	P=4	<b>1.403</b>	<b>1.315</b>	<b>1.423</b>	<b>1.282</b>

**More Results on RNN** Table 1 shows the final NLL for various setup of RNNs and algorithms. Lower NLLs can be achieved when using 4 workers compared to 1 worker, *i.e.*, standard SG-MCMC.

**More Results  $t$ -SNE Visualization** Following the same procedure of the  $t$ -SNE visualization for SVHN, we also embed the 10,000 test digits of MNIST in Fig. 3, using 20 model samples of the 400-400 FNN, obtained by downpour pSGLD algorithm. Though many methods can achieve high classification accuracy on MNIST dataset, there are some digits that are intrinsically hard to classify. In practice, it may be beneficial to leverage the uncertainty to further boost performance.

## References

- C. Chen, N. Ding, and L. Carin. On the convergence of stochastic gradient MCMC algorithms with high-order integrators. In *NIPS*, 2015.
- C. Chen, N. Ding, C. Li, Y. Zhang, and L. Carin. Stochastic gradient MCMC with stale gradients. In *NIPS*, 2016.
- C. Li, C. Chen, D. Carlson, and L. Carin. Preconditioned stochastic gradient Langevin dynamics for deep neural networks. In *AAAI*, 2016.
- Y. W. Teh, A. H. Thiéry, and S. J. Vollmer. Consistency and fluctuations for stochastic gradient Langevin dynamics. *JMLR*, 2016.
- S. J. Vollmer, K. C. Zygalakis, and Y. W. Teh. (Non-)asymptotic properties of stochastic gradient Langevin dynamics. Technical report, 2015.