

## SUPPLEMENTARY MATERIAL

Our unified analysis of ASYNC-RED is based on the monotone operator theory (Ryu & Boyd, 2016). In Supplement A, we first clarify our setting for the access of the shared memory. In Supplement B, we present the proof of Theorem 1 and Theorem 2, proving the fixed-point convergence of ASYNC-RED to  $\text{zer}(\mathbf{G})$  in both batch and stochastic settings. In Supplement C, we provide a brief review of the related knowledge on monotone operators. In Supplement D, we include additional technical details and experiments omitted from the main paper due to space.

### A MEMORY ACCESS WITHOUT GLOBAL LOCK

In the setting of ASYNC-RED, multiple cores may simultaneously read and update the blocks  $\mathbf{x}_i$  in shared memory. We coordinate the memory access of different cores by imposing certain *local* locks. For example, consider one work cycle of core  $c_i$  for updating the block  $\mathbf{x}_i$ . First, a local *read* lock is imposed to  $\mathbf{x}_i$  such that only read operations (by  $c_i$  or others) can be performed on  $\mathbf{x}_i$ . If, at the same time, other cores want to write  $\mathbf{x}_i$ , then they have to wait until the read lock is released by the last one who finishes reading the block. However, if they want to write other blocks, their operations will not be blocked. Secondly, core  $c_i$  evaluates the RED update on  $\mathbf{x}_i$ , while other cores continuously update  $\mathbf{x}$ . Here, we assume that the number of updates by cores other than  $c_i$  is bounded by some positive integer, which is exactly what Assumption 1 refers to. After the evaluation finishes, core  $c_i$  imposes a local *write* lock, which prevents both read and write by other cores, on  $\mathbf{x}_i$  and write the block with the computed update. Similarly, other cores have to wait until the lock is released before operating on  $\mathbf{x}_i$ . Finally, when the update finishes, the local lock will be released and core  $c_i$  will restart a new cycle. Note that  $\mathbf{x}$  is never locked *globally* during the full update cycle, and the reads of each block are always consistent.

In order to ensure the consistent read of  $\mathbf{x}$ , we leverage the dual-memory strategy for block coordinate settings proposed in (Peng et al., 2016) (see section 1.2.1 ‘Block coordinate’). Its key idea is that, before every write to a block  $\mathbf{x}_i$ , a copy of the old version of the block is kept for reading. In this way, there always exists some state of  $\mathbf{x}$  in the memory for the cores to access.

### B PROOF OF ANALYSIS

In this section, we first present the proof of Theorem 1, then followed by the proof of Theorem 2. For a review of monotone operators, we refer to Supplement C.

Throughout the proof, we consider the probability space  $(\Omega, \mathcal{F}, P)$ , where  $\Omega$  denotes the sample space,  $\mathcal{F}$  the  $\sigma$ -algebra, and  $P$  the probability measure.  $\mathbf{x}^k$  is a random variable defined in  $\mathbb{R}^n$ . We use  $\|\cdot\|$  to denote the  $\ell_2$ -norm. We define the sequence of sub  $\sigma$ -algebra  $\{\mathcal{X}^k\}_{k \in \mathbb{N}}$  of  $\mathcal{F}$  as

$$\mathcal{X}^k := \sigma(\mathbf{x}^0, \dots, \mathbf{x}^k, \Delta_0, \dots, \Delta_k),$$

where  $\sigma$  generates the filtration (smallest  $\sigma$ -algebra) from  $\mathbf{x}^0, \dots, \mathbf{x}^k$ , and  $\Delta_0, \dots, \Delta_k$ . Note that the sequence  $\{\mathcal{X}^k\}_{k \in \mathbb{N}}$  is such that  $\mathcal{X}^k \subset \mathcal{X}^{k+1}$  for any  $k \in \mathbb{N}$ . We use  $\mathbf{x}^*$  to denote some fixed point in the set  $\text{zer}(\mathbf{G})$ .

#### B.1 PROOF OF THEOREM 1

Our proof needs the following lemma on the RED operator.

**Lemma 1.** *Let Assumption 3 and 4 hold for  $g$  and  $\mathbf{D}_\sigma$ . The composite operator  $\mathbf{G}$  is  $1/(L + 2\tau)$ -cocoercive, that is*

$$(\mathbf{G}(\mathbf{x}) - \mathbf{G}(\mathbf{y}))^\top (\mathbf{x} - \mathbf{y}) \geq \frac{1}{L + 2\tau} \|\mathbf{G}(\mathbf{x}) - \mathbf{G}(\mathbf{y})\|^2.$$

*Proof.* This lemma is adapted from Lemma 3 in (Sun et al., 2019a). Consider the following decomposition

$$\mathbf{I} - \frac{2}{L + 2\tau} \mathbf{G} = \left( \frac{2}{L + 2\tau} \cdot \frac{L}{2} \right) \left[ \mathbf{I} - \frac{2}{L} \nabla g \right] + \left( \frac{2}{L + 2\tau} \cdot \frac{2\tau}{2} \right) \left[ \mathbf{I} - \frac{1}{\tau} \mathbf{H} \right], \quad (16)$$

where we recall  $H = \tau(I - D_\sigma)$ . According to Assumption [3](#),  $g$  is convex and  $\nabla g$  is  $L$ -Lipschitz continuous. By Proposition [1](#) in Supplement [C](#),  $\nabla g$  is  $1/L$ -cocoercive. Hence, by Proposition [2](#) in Supplement [C](#),  $I - (2/L)\nabla g$  is nonexpansive. Since  $D_\sigma = I - (1/\tau)H$ , this means that  $I - (1/\tau)H$  is nonexpansive. From Proposition [3](#) in Supplement [C](#), we know that the convex combination of two nonexpansive operators is nonexpansive. Thus,  $I - (2/(L + 2\tau))G$  is nonexpansive, which also means that  $G$  is  $1/(L + 2\tau)$ -cocoercive according to Proposition [2](#) in Supplement [C](#).  $\square$

Now we can start the main proof. Under the fixed stepsize  $\gamma > 0$ , we begin with the following equations regarding the fixed point  $\mathbf{x}^* \in \text{zer}(G)$

$$\begin{aligned} & \mathbb{E} [\|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 | \mathcal{X}^k] \\ &= \mathbb{E} [\|\mathbf{x}^k - \gamma G_i(\tilde{\mathbf{x}}^k) - \mathbf{x}^*\|^2 | \mathcal{X}^k] \\ &= \mathbb{E} [\|\mathbf{x}^k - \mathbf{x}^*\|^2 | \mathcal{X}^k] + \gamma^2 \mathbb{E} [\|G_i(\tilde{\mathbf{x}}^k)\|^2 | \mathcal{X}^k] + 2\gamma \mathbb{E} [(G_i(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) | \mathcal{X}^k] \end{aligned} \quad (17)$$

Since  $G_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is evaluated on a random block of  $\mathbf{x}_i$ , we have the following conditional expectations

$$\mathbb{E} [(G_i(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) | \mathcal{X}^k] = \frac{1}{b} \sum_{i=1}^b (G_i(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) = \frac{1}{b} (G(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \quad (18)$$

and

$$\mathbb{E} [\|G_i(\tilde{\mathbf{x}}^k)\|^2 | \mathcal{X}^k] = \frac{1}{b} \sum_{i=1}^b \|G_i(\tilde{\mathbf{x}}^k)\|^2 = \frac{1}{b} \|G(\tilde{\mathbf{x}}^k)\|^2. \quad (19)$$

Thus, plugging the above results into [\(17\)](#)

$$\mathbb{E} [\|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 | \mathcal{X}^k] \leq \|\mathbf{x}^k - \mathbf{x}^*\|^2 + \frac{\gamma^2}{b} \|G(\tilde{\mathbf{x}}^k)\|^2 + \underbrace{\frac{2\gamma}{b} (G(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k)}_{(\dagger)}. \quad (20)$$

The term  $(\dagger)$  can be expressed as

$$\begin{aligned} & \frac{2\gamma}{b} (G(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \\ &= \frac{2\gamma}{b} (G(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \tilde{\mathbf{x}}^k + \sum_{s=k-\Delta_k}^{k-1} (\mathbf{x}^s - \mathbf{x}^{s+1})) \\ &= \frac{2\gamma}{b} (G(\tilde{\mathbf{x}}^k) - G(\mathbf{x}^*))^\top (\mathbf{x}^* - \tilde{\mathbf{x}}^k) + \frac{2\gamma}{b} (G(\tilde{\mathbf{x}}^k))^\top \left( \sum_{s=k-\Delta_k}^{k-1} (\mathbf{x}^s - \mathbf{x}^{s+1}) \right) \\ &= \frac{2\gamma}{b} (G(\tilde{\mathbf{x}}^k) - G(\mathbf{x}^*))^\top (\mathbf{x}^* - \tilde{\mathbf{x}}^k) + \frac{2\gamma^2}{b} \sum_{s=k-\Delta_k}^{k-1} G(\tilde{\mathbf{x}}^k)^\top G_{i_s}(\tilde{\mathbf{x}}^s), \end{aligned} \quad (21)$$

where in the second line we used the definition of the stale iterate  $\mathbf{x}^{s+1} = \mathbf{x}^s - \gamma G_{i_s}(\tilde{\mathbf{x}}^k)$ , and in the third line the fact that  $G(\mathbf{x}^*) = \mathbf{0}$ . By using Lemma [1](#), we obtain the upper bound for the first term in equation [\(21\)](#)

$$\frac{2\gamma}{b} (G(\tilde{\mathbf{x}}^k) - G(\mathbf{x}^*))^\top (\mathbf{x}^* - \tilde{\mathbf{x}}^k) \leq -\frac{2\gamma \|G(\tilde{\mathbf{x}}^k)\|^2}{b(L + 2\tau)}. \quad (22)$$

For the second term in [\(21\)](#), we have

$$\begin{aligned} \frac{2\gamma^2}{b} \sum_{s=k-\Delta_k}^{k-1} G(\tilde{\mathbf{x}}^k)^\top G_{i_s}(\tilde{\mathbf{x}}^s) &\leq \frac{\lambda\gamma^2 \|G(\tilde{\mathbf{x}}^k)\|^2}{b} + \sum_{s=k-\Delta_k}^{k-1} \frac{\gamma^2 \|G_{i_s}(\tilde{\mathbf{x}}^s)\|^2}{b}, \\ &\leq \frac{\lambda\gamma^2 \|G(\tilde{\mathbf{x}}^k)\|^2}{b} + \sum_{s=k-\lambda}^{k-1} \frac{\gamma^2 \|G(\tilde{\mathbf{x}}^s)\|^2}{b}, \end{aligned} \quad (23)$$

where in the first inequality we used the *Young's inequality*

$$\mathbf{x}_1^\top \mathbf{x}_2 \leq \frac{1}{2} [\|\mathbf{x}_1\|^2 + \|\mathbf{x}_2\|^2], \quad (24)$$

and in the second inequality we use

$$\sum_{s=k-\Delta k}^{k-1} \gamma^2 \|\mathbf{G}_{i_s}(\tilde{\mathbf{x}}^s)\|^2 = \sum_{s=k-\Delta k}^{k-1} \|\mathbf{x}^s - \mathbf{x}^{s+1}\|_2^2 \leq \sum_{s=k-\lambda}^{k-1} \|\mathbf{x}^s - \mathbf{x}^{s+1}\|_2^2 = \sum_{s=k-\lambda}^{k-1} \gamma^2 \|\mathbf{G}(\tilde{\mathbf{x}}^s)\|^2.$$

Applying (22) and (23) in (21) yields the overall upper bound for the term (†)

$$\frac{2\gamma}{b} (\mathbf{G}(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \leq \frac{(L+2\tau)\lambda\gamma^2 - 2\gamma}{(L+2\tau)b} \|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2 + \sum_{s=k-\lambda}^{k-1} \frac{\gamma^2 \|\mathbf{G}(\tilde{\mathbf{x}}^s)\|^2}{b}. \quad (25)$$

Next, by plugging (25) into (17) and re-arranging the terms, we obtain the following inequality

$$\begin{aligned} & \mathbb{E} [\|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 | \mathcal{X}^k] \\ & \leq \|\mathbf{x}^k - \mathbf{x}^*\|^2 + \sum_{s=k-\lambda}^{k-1} \frac{\gamma^2 \|\mathbf{G}(\tilde{\mathbf{x}}^s)\|^2}{b} + \frac{(L+2\tau)(1+\lambda)\gamma^2 - 2\gamma}{(L+2\tau)b} \|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2. \end{aligned} \quad (26)$$

Taking the total expectation of equation (26) and re-arranging the terms yields that

$$\begin{aligned} & \frac{2\gamma - (L+2\tau)(1+\lambda)\gamma^2}{(L+2\tau)b} \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \\ & \leq \mathbb{E} [\|\mathbf{x}^k - \mathbf{x}^*\|^2] - \mathbb{E} [\|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2] + \gamma^2 \sum_{s=k-\lambda}^{k-1} \frac{\mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^s)\|^2]}{b} \end{aligned} \quad (27)$$

We then telescope-sum equation (27) over  $t > 0$  iterations to have

$$\begin{aligned} & \sum_{k=0}^{t-1} \frac{2\gamma - (L+2\tau)(1+\lambda)\gamma^2}{(L+2\tau)b} \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \\ & \leq \mathbb{E} [\|\mathbf{x}^0 - \mathbf{x}^*\|^2] - \mathbb{E} [\|\mathbf{x}^t - \mathbf{x}^*\|^2] + \gamma^2 \sum_{k=0}^{t-1} \sum_{s=k-\lambda}^{k-1} \frac{\mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^s)\|^2]}{b} \end{aligned} \quad (28)$$

where the index  $s$  always start at 0. Under the assumption of consistent read, it is true that

$$\sum_{k=0}^{t-1} \sum_{s=k-\lambda}^{k-1} \frac{\mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^s)\|^2]}{b} \leq \lambda \sum_{k=0}^{t-1} \frac{\mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2]}{b}. \quad (29)$$

In the case of inconsistent read, the above inequality does not always hold. We refer to Peng et al. (2016) for a comprehensive analysis for asynchronous block-coordinate methods with inconsistent reads. Now, we rewrite equation (28) as

$$\sum_{k=0}^{t-1} \frac{2\gamma - (L+2\tau)(1+2\lambda)\gamma^2}{(L+2\tau)b} \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \leq \mathbb{E} [\|\mathbf{x}^0 - \mathbf{x}^*\|^2] - \mathbb{E} [\|\mathbf{x}^t - \mathbf{x}^*\|^2]. \quad (30)$$

In order to ensure the convergence, we need the coefficient of  $\mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2]$  to be positive. From basic algebra, one feasible range for the stepsize  $\gamma$  is

$$0 < \gamma \leq \frac{1}{(L+2\tau)(1+2\lambda)},$$

which directly implies that

$$0 < \frac{\gamma}{(L+2\tau)b} \leq \frac{2\gamma - (L+2\tau)(1+2\lambda)\gamma^2}{(L+2\tau)b}.$$

By simplifying (30) with the above result and dropping the negative term, we can derive the following bound for the  $\mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2]$  averaged over  $t$  iterations

$$\frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \leq \frac{(L+2\tau)b}{\gamma t} \mathbb{E} [\|\mathbf{x}^0 - \mathbf{x}^*\|^2] \leq \frac{(L+2\tau)b}{\gamma t} R_0^2. \quad (31)$$

The above inequality establishes that the change of the stale iterate  $\tilde{\mathbf{x}}^k$  converges to zero as  $t$  increases. Next, we will use the bound to establish the similar result for the actual iterate  $\mathbf{x}^k$ . We know that  $\|\mathbf{G}(\mathbf{x}^k)\|^2$  can be bounded by

$$\begin{aligned} \|\mathbf{G}(\mathbf{x}^k)\|^2 &\leq (\|\mathbf{G}(\mathbf{x}^k) - \mathbf{G}(\tilde{\mathbf{x}}^k)\| + \|\mathbf{G}(\tilde{\mathbf{x}}^k)\|)^2 \\ &= \|\mathbf{G}(\mathbf{x}^k) - \mathbf{G}(\tilde{\mathbf{x}}^k)\|^2 + \|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2 + 2\|\mathbf{G}(\mathbf{x}^k) - \mathbf{G}(\tilde{\mathbf{x}}^k)\| \|\mathbf{G}(\tilde{\mathbf{x}}^k)\| \\ &\leq 2\|\mathbf{G}(\mathbf{x}^k) - \mathbf{G}(\tilde{\mathbf{x}}^k)\|^2 + 2\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2 \\ &\leq 2(L+2\tau)^2 \|\mathbf{x}^k - \tilde{\mathbf{x}}^k\|^2 + 2\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2 \end{aligned} \quad (32)$$

where in the second inequality we used the Young's inequality (24), and in the third inequality we used the following result implied by Lemma 1

$$(L+2\tau)\|\mathbf{x} - \mathbf{y}\| \geq \|\mathbf{G}(\mathbf{x}) - \mathbf{G}(\mathbf{y})\|.$$

By expressing the stale iterate  $\tilde{\mathbf{x}}^k$ , we can write equation (32) as

$$\begin{aligned} \|\mathbf{G}(\mathbf{x}^k)\|^2 &\leq 2(L+2\tau)^2 \left\| \sum_{s=k-\lambda}^{k-1} \gamma \mathbf{G}_{i_s}(\tilde{\mathbf{x}}^s) \right\|^2 + 2\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2 \\ &\leq 2\lambda(L+2\tau)^2 \sum_{s=k-\lambda}^{k-1} \gamma^2 \|\mathbf{G}_{i_s}(\tilde{\mathbf{x}}^s)\|^2 + 2\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2. \end{aligned} \quad (33)$$

where we use the fact

$$\left\| \sum_{i=1}^n \mathbf{x}_i \right\|^2 = \sum_{i=1}^n \|\mathbf{x}_i\|^2 + \sum_{a \neq b} \mathbf{x}_a^\top \mathbf{x}_b \leq \sum_{i=1}^n \|\mathbf{x}_i\|^2 + \frac{1}{2} \sum_{a \neq b} [\|\mathbf{x}_a\|^2 + \|\mathbf{x}_b\|^2] = n \sum_{i=1}^n \|\mathbf{x}_i\|^2$$

Taking the expectation of equation (33) leads to

$$\begin{aligned} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] &\leq 2\lambda(L+2\tau)^2 \sum_{s=k-\lambda}^{k-1} \gamma^2 \mathbb{E} [\|\mathbf{G}_{i_s}(\tilde{\mathbf{x}}^s)\|^2] + 2\mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \\ &\leq 2\lambda(L+2\tau)^2 \sum_{s=k-\lambda}^{k-1} \frac{\gamma^2 \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^s)\|^2]}{b} + 2\mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2], \end{aligned} \quad (34)$$

By averaging (34) over  $t > 0$  iterations, we obtain that

$$\begin{aligned} \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] &\leq \frac{2\lambda(L+2\tau)^2}{t} \sum_{k=0}^{t-1} \sum_{s=k-\lambda}^{k-1} \frac{\gamma^2 \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^s)\|^2]}{b} + \frac{2}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \\ &\leq \frac{2\lambda^2(L+2\tau)^2}{t} \sum_{k=0}^{t-1} \frac{\gamma^2 \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2]}{b} + \frac{2}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \end{aligned} \quad (35)$$

where we again used result in (29) in the last inequality. Re-arranging the terms in (35) yields

$$\frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \left[ \frac{2\lambda^2(L+2\tau)^2}{b} \gamma^2 + 2 \right] \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \quad (36)$$

We plug the result in (31) into (36) and obtain

$$\frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \left[ \frac{2\lambda^2(L+2\tau)^2}{b} \gamma^2 + 2 \right] \frac{(L+2\tau)b}{\gamma t} R_0^2, \quad (37)$$

Since it is always true that

$$\gamma \leq \frac{1}{(L+2\tau)(1+2\lambda)} \leq \frac{1}{(L+2\tau)(1+\lambda)}.$$

we can simplify the bound by using the above inequality related to the stepsize  $\gamma$

$$\frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \left[ \frac{2\lambda^2}{(1+\lambda)^2 b} + 2 \right] \frac{(L+2\tau)b}{\gamma t} R_0^2. \quad (38)$$

Let  $D = 2\lambda^2/(1+\lambda)^2$ , and we derive the desired result.

$$\min_{0 \leq k \leq t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \left[ \frac{D}{b} + 2 \right] \frac{(L+2\tau)b}{\gamma t} R_0^2. \quad (39)$$

## B.2 PROOF OF THEOREM 2

We prove Theorem 2 by following the procedure in the proof of Theorem 1 with the adaptation to the block stochastic operator  $\widehat{\mathbf{G}}_i$ . In the key steps, we will highlight the difference between the two proofs. In addition to Lemma 1, our second proof requires the following lemma related to the statistical properties of  $\widehat{\mathbf{G}}$ .

**Lemma 2.** *Let Assumption 3 and 4 hold for  $g$  and  $\mathbf{D}_\sigma$ . Then, we can establish the following statements for operator  $\widehat{\mathbf{G}}$*

$$\mathbb{E} [\widehat{\mathbf{G}}(\mathbf{x})] = \mathbf{G}(\mathbf{x}), \quad \mathbb{E} [\|\widehat{\mathbf{G}}(\mathbf{x}) - \mathbf{G}(\mathbf{x})\|^2] \leq \frac{\nu^2}{w},$$

which further implies that

$$\mathbb{E} [\|\widehat{\mathbf{G}}(\mathbf{x})\|^2] \leq \frac{\nu^2}{w} + \|\mathbf{G}(\mathbf{x})\|^2.$$

*Proof.* Since the the stochasticity happens only in the evaluation of the gradient, it is straightforward to see that

$$\mathbb{E} [\widehat{\mathbf{G}}(\mathbf{x})] = \mathbb{E} [\widehat{\nabla} g(\mathbf{x})] + \mathbf{D}_\sigma(\mathbf{x}) = \mathbf{G}(\mathbf{x}),$$

Similarly, we have that

$$\mathbb{E} [\|\widehat{\mathbf{G}}(\mathbf{x}) - \mathbf{G}(\mathbf{x})\|_2^2] = \mathbb{E} [\|\widehat{\nabla} g(\mathbf{x}) - \nabla g(\mathbf{x})\|_2^2] \leq \frac{\nu^2}{w}$$

Given that  $\text{Tr}(\mathbb{E} [X^\top X]) = \text{Tr}(\text{Cov} [X]) + \text{Tr}(\mathbb{E} [X]^2)$ , we obtain that

$$\mathbb{E} [\|\widehat{\mathbf{G}}(\mathbf{x})\|^2] = \mathbb{E} [\|\widehat{\mathbf{G}}(\mathbf{x}) - \mathbf{G}(\mathbf{x})\|^2] + \mathbb{E} [\|\mathbf{G}(\mathbf{x})\|^2] \leq \frac{\nu^2}{w} + \|\mathbf{G}(\mathbf{x})\|^2,$$

where we let  $\mathbb{E} [\widehat{\mathbf{G}}(\mathbf{x})]^2 := \mathbb{E} [\widehat{\mathbf{G}}(\mathbf{x})]^\top \mathbb{E} [\widehat{\mathbf{G}}(\mathbf{x})]$ . Note that  $\text{Tr}(\cdot)$  and  $\text{Cov}(\cdot)$  denote the computation of the trace and covariance of a matrix and a vector, respectively.  $\square$

Now we start the proof. Similar as (17), we write that

$$\begin{aligned} & \mathbb{E} [\|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 | \mathcal{X}^k] \\ &= \mathbb{E} [\|\mathbf{x}^k - \gamma \widehat{\mathbf{G}}_i(\tilde{\mathbf{x}}^k) - \mathbf{x}^*\|^2 | \mathcal{X}^k] \\ &= \mathbb{E} [\|\mathbf{x}^k - \mathbf{x}^*\|^2 | \mathcal{X}^k] + \gamma^2 \mathbb{E} [\|\widehat{\mathbf{G}}_i(\tilde{\mathbf{x}}^k)\|^2 | \mathcal{X}^k] + 2\gamma \mathbb{E} [(\widehat{\mathbf{G}}_i(\tilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) | \mathcal{X}^k] \end{aligned} \quad (40)$$

Here, the conditional expectation is taken for  $\widehat{G}_i(\mathbf{x}) = \mathbf{U}_i \mathbf{U}_i^\top \widehat{G}(\mathbf{x})$ . By using Lemma 2, we can compute conditional expectations as

$$\mathbb{E} \left[ (\widehat{G}_i(\widetilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) | \mathcal{X}^k \right] = \frac{1}{b} \mathbb{E} \left[ (\widehat{G}(\widetilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) | \mathcal{X}^k \right] = \frac{1}{b} (G(\widetilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \quad (41)$$

and

$$\mathbb{E} \left[ \|\widehat{G}_i(\widetilde{\mathbf{x}}^k)\|^2 | \mathcal{X}^k \right] = \frac{1}{b} \mathbb{E} \left[ \|\widehat{G}(\widetilde{\mathbf{x}}^k)\|^2 | \mathcal{X}^k \right] \leq \frac{\nu^2}{wb} + \frac{\|G(\widetilde{\mathbf{x}}^k)\|^2}{b}. \quad (42)$$

where we first compute the expectation corresponding to the randomized block and then the expectation for the stochastic measurements. We note that the expectation of the cross term (41) remains the same as the result in (18), while the expectation in (42) has one extra term related to the norm variance of the stochastic operator compared with (19). As we shall see in the future steps, the difference in the expectation of the operator's squared norm leads to the most modifications. Using the above results in equation (40) yields that

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 | \mathcal{X}^k \right] \\ & \leq \|\mathbf{x}^k - \mathbf{x}^*\|^2 + \frac{\gamma^2}{b} \|G(\widetilde{\mathbf{x}}^k)\|^2 + \frac{\gamma^2 \nu^2}{wb} + \underbrace{\frac{2\gamma}{b} (G(\widetilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k)}_{(\dagger)}. \end{aligned} \quad (43)$$

By following (21), we can express the term  $(\dagger)$  as

$$\begin{aligned} & \frac{2\gamma}{b} (G(\widetilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \\ & = \frac{2\gamma}{b} (G(\widetilde{\mathbf{x}}^k) - G(\mathbf{x}^*))^\top (\mathbf{x}^* - \widetilde{\mathbf{x}}^k) + \frac{2\gamma^2}{b} \sum_{s=k-\Delta_k}^{k-1} G(\widetilde{\mathbf{x}}^k)^\top \widehat{G}_{i_s}(\widetilde{\mathbf{x}}^s), \end{aligned} \quad (44)$$

The upper bound of the first term is the same as shown in (22), which is

$$\frac{2\gamma}{b} (G(\widetilde{\mathbf{x}}^k) - G(\mathbf{x}^*))^\top (\mathbf{x}^* - \widetilde{\mathbf{x}}^k) \leq -\frac{2\gamma \|G(\widetilde{\mathbf{x}}^k)\|^2}{b(L+2\tau)}. \quad (45)$$

Similarly, our second term is bounded by

$$\frac{2\gamma^2}{b} \sum_{s=k-\Delta_k}^{k-1} G(\widetilde{\mathbf{x}}^k)^\top \widehat{G}_{i_s}(\widetilde{\mathbf{x}}^s) \leq \frac{\lambda \gamma^2 \|G(\widetilde{\mathbf{x}}^k)\|^2}{b} + \sum_{s=k-\lambda}^{k-1} \frac{\gamma^2 \|\widehat{G}(\widetilde{\mathbf{x}}^s)\|^2}{b}, \quad (46)$$

where we used the Young's inequality (24) together with the fact that

$$\sum_{s=k-\Delta_k}^{k-1} \|\widehat{G}_{i_s}(\widetilde{\mathbf{x}}^k)\|^2 \leq \sum_{s=k-\lambda}^{k-1} \|\widehat{G}_{i_s}(\widetilde{\mathbf{x}}^k)\|^2 \leq \sum_{s=k-\lambda}^{k-1} \|\widehat{G}(\widetilde{\mathbf{x}}^k)\|^2.$$

Equation (45) and (46) together establish the overall upper bound for the term  $(\dagger)$

$$\frac{2\gamma}{b} (G(\widetilde{\mathbf{x}}^k))^\top (\mathbf{x}^* - \mathbf{x}^k) \leq \frac{(L+2\tau)\lambda\gamma^2 - 2\gamma}{(L+2\tau)b} \|G(\widetilde{\mathbf{x}}^k)\|^2 + \sum_{s=k-\lambda}^{k-1} \frac{\gamma^2 \|\widehat{G}(\widetilde{\mathbf{x}}^s)\|^2}{b}. \quad (47)$$

By plugging (47) into (40) and re-arranging the terms, we obtain that

$$\begin{aligned} & \mathbb{E} \left[ \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 | \mathcal{X}^k \right] \\ & \leq \|\mathbf{x}^k - \mathbf{x}^*\|^2 + \frac{\gamma^2 \nu^2}{wb} + \sum_{s=k-\lambda}^{k-1} \frac{\gamma^2 \|\widehat{G}(\widetilde{\mathbf{x}}^s)\|^2}{b} + \frac{(L+2\tau)(1+\lambda)\gamma^2 - 2\gamma}{(L+2\tau)b} \|G(\widetilde{\mathbf{x}}^k)\|^2. \end{aligned} \quad (48)$$

Taking the total expectation of equation (48) and re-arranging the terms yields that

$$\begin{aligned} & \frac{2\gamma - (L+2\tau)(1+\lambda)\gamma^2}{(L+2\tau)b} \mathbb{E} \left[ \|G(\widetilde{\mathbf{x}}^k)\|^2 \right] \\ & \leq \mathbb{E} \left[ \|\mathbf{x}^k - \mathbf{x}^*\|^2 \right] - \mathbb{E} \left[ \|\mathbf{x}^{k+1} - \mathbf{x}^*\|^2 \right] + \frac{\gamma^2 \nu^2}{wb} + \gamma^2 \sum_{s=k-\lambda}^{k-1} \left[ \frac{\nu^2}{wb} + \frac{\mathbb{E} \left[ \|G(\widetilde{\mathbf{x}}^s)\|^2 \right]}{b} \right] \end{aligned} \quad (49)$$

where we use the following inequality derived by using the law of total expectation and Lemma 2

$$\mathbb{E} [\|\widehat{\mathbf{G}}(\widetilde{\mathbf{x}}^s)\|^2] = \mathbb{E} [\mathbb{E} [\|\widehat{\mathbf{G}}(\widetilde{\mathbf{x}}^s)\|^2 | \mathcal{X}^s]] \leq \frac{\nu^2}{w} + \mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^s)\|^2]. \quad (50)$$

We telescope-sum equation (49) over  $t > 0$  iterations to obtain

$$\begin{aligned} & \sum_{k=0}^{t-1} \frac{2\gamma - (L + 2\tau)(1 + \lambda)\gamma^2}{(L + 2\tau)b} \mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2] \\ & \leq \mathbb{E} [\|\mathbf{x}^0 - \mathbf{x}^*\|^2] - \mathbb{E} [\|\mathbf{x}^t - \mathbf{x}^*\|^2] + \sum_{k=0}^{t-1} \frac{\gamma^2 \nu^2}{wb} + \gamma^2 \sum_{k=0}^{t-1} \sum_{s=k-\lambda}^{k-1} \left[ \frac{\nu^2}{wb} + \frac{\mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^s)\|^2]}{b} \right] \end{aligned} \quad (51)$$

By applying the same relaxation trick in (29) to (51)

$$\sum_{k=0}^{t-1} \sum_{s=k-\lambda}^{k-1} \left[ \frac{\nu^2}{wb} + \frac{\mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^s)\|^2]}{b} \right] \leq \lambda \sum_{k=0}^{t-1} \left[ \frac{\nu^2}{wb} + \frac{\mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2]}{b} \right], \quad (52)$$

we then have that

$$\sum_{k=0}^{t-1} \frac{2\gamma - (L + 2\tau)(1 + 2\lambda)\gamma^2}{(L + 2\tau)b} \mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2] \leq \mathbb{E} [\|\mathbf{x}^0 - \mathbf{x}^*\|^2] + \frac{(1 + \lambda)\gamma^2 \nu^2}{wb} \cdot t, \quad (53)$$

where we dropped the negative term. Recall that if  $\gamma$  is in the range  $\gamma \in (0, 1/((L + 2\tau)(1 + 2\lambda))]$ , we have the inequality

$$\frac{\gamma}{(L + 2\tau)b} \leq \frac{2\gamma - (L + 2\tau)(1 + 2\lambda)\gamma^2}{(L + 2\tau)b}.$$

By relaxing the coefficient in the lefthand side, dividing the inequality by  $t$ , and re-arranging the terms, we obtain the convergence in terms of the stale iterate  $\widetilde{\mathbf{x}}^k$

$$\begin{aligned} \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2] & \leq \frac{(L + 2\tau)b}{\gamma t} \left[ \mathbb{E} [\|\mathbf{x}^0 - \mathbf{x}^*\|^2] + \frac{(1 + \lambda)\gamma^2 \nu^2}{wb} \cdot t \right] \\ & \leq \frac{(L + 2\tau)b}{\gamma t} R_0^2 + \frac{\gamma}{w} C \end{aligned} \quad (54)$$

where we used Assumption 2 and let  $C = (L + 2\tau)(1 + \lambda)\nu^2$ . Compared with the result in equation (31), equation (54) has the extra term related to the variance of  $\widehat{\mathbf{G}}_i(\mathbf{x})$ . Next, we establish the convergence in terms of actual iterate  $\mathbf{x}^k$ . Following the steps from (32) to (34), we directly obtain the inequality related to  $\widehat{\mathbf{G}}_i(\widetilde{\mathbf{x}})$

$$\mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq 2\lambda(L + 2\tau)^2 \sum_{s=k-\lambda}^{k-1} \gamma^2 \mathbb{E} [\|\widehat{\mathbf{G}}_{i_s}(\widetilde{\mathbf{x}}^s)\|^2] + 2\mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2] \quad (55)$$

By using the the result in (50), we derive from (55) that

$$\mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq 2\lambda(L + 2\tau)^2 \sum_{s=k-\lambda}^{k-1} \gamma^2 \left[ \frac{\nu^2}{wb} + \frac{\mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^s)\|^2]}{b} \right] + 2\mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2]. \quad (56)$$

By averaging (56) over  $t > 0$  iterations, we obtain that

$$\begin{aligned} & \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \\ & \leq \frac{2\lambda(L + 2\tau)^2}{t} \sum_{k=0}^{t-1} \sum_{s=k-\lambda}^{k-1} \gamma^2 \left[ \frac{\nu^2}{wb} + \frac{\mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^s)\|^2]}{b} \right] + \frac{2}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2] \\ & \leq \frac{2\lambda^2(L + 2\tau)^2}{t} \sum_{k=0}^{t-1} \gamma^2 \left[ \frac{\nu^2}{wb} + \frac{\mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2]}{b} \right] + \frac{2}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\widetilde{\mathbf{x}}^k)\|^2] \end{aligned} \quad (57)$$

where we again used the relaxation (52) in the last inequality. Re-arranging the terms in (57) yields

$$\begin{aligned} & \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \\ & \leq \frac{2\lambda^2(L+2\tau)^2 \cdot \nu^2}{wb} \gamma^2 + \left[ \frac{2\lambda^2(L+2\tau)^2}{b} \gamma^2 + 2 \right] \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\tilde{\mathbf{x}}^k)\|^2] \end{aligned} \quad (58)$$

We plug the result in (54) into (58) and obtain

$$\begin{aligned} & \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \\ & \leq \frac{2\lambda^2(L+2\tau)^2 \cdot \nu^2}{wb} \gamma^2 + \left[ \frac{2\lambda^2(L+2\tau)^2}{b} \gamma^2 + 2 \right] \left[ \frac{(L+2\tau)b}{\gamma t} R_0^2 + \frac{\gamma}{w} C \right] \end{aligned} \quad (59)$$

Similarly, we can use the fact

$$\gamma \leq \frac{1}{(L+2\tau)(1+\lambda)}.$$

to simplify the bound in (59)

$$\begin{aligned} & \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \\ & \leq \frac{2\lambda^2(L+2\tau)^2 \cdot \nu^2}{wb} \cdot \frac{1}{(L+2\tau)(1+\lambda)} \cdot \gamma + \left[ \frac{2\lambda^2}{(1+\lambda)^2 b} + 2 \right] \left[ \frac{(L+2\tau)b}{\gamma t} R_0^2 + \frac{\gamma}{w} C \right] \\ & = \frac{2\lambda^2}{(1+\lambda)^2 b} \cdot \frac{(L+2\tau)(1+\lambda)\nu^2}{w} \cdot \gamma + \left[ \frac{2\lambda^2}{(1+\lambda)^2 b} + 2 \right] \left[ \frac{(L+2\tau)b}{\gamma t} R_0^2 + \frac{\gamma}{w} C \right] \\ & = \frac{2\lambda^2}{(1+\lambda)^2 b} \cdot \frac{C}{w} \gamma + \left[ \frac{2\lambda^2}{(1+\lambda)^2 b} + 2 \right] \left[ \frac{(L+2\tau)b}{\gamma t} R_0^2 + \frac{\gamma}{w} C \right] \end{aligned} \quad (60)$$

where we recall  $C = (L+2\tau)(1+\lambda)\nu^2$ . Let  $D = 2\lambda^2/(1+\lambda)^2$  and we can derive the result of Theorem 2

$$\min_{0 \leq k \leq t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \left[ \frac{D}{b} + 2 \right] \frac{(L+2\tau)b}{\gamma t} R_0^2 + \left[ \frac{2D}{b} + 2 \right] \frac{\gamma}{w} C, \quad (61)$$

which immediately implies the result in remark 1 by setting  $\gamma = 1/\sqrt{wt}$

$$\min_{0 \leq k \leq t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \frac{1}{t} \sum_{k=0}^{t-1} \mathbb{E} [\|\mathbf{G}(\mathbf{x}^k)\|^2] \leq \left[ \frac{D}{b} + 2 \right] \frac{(L+2\tau)b}{\sqrt{wt}} R_0^2 + \left[ \frac{2D}{b} + 2 \right] \frac{C}{\sqrt{wt}}. \quad (62)$$

From basic algebra, we can derive the condition for  $\lambda$

$$\frac{1}{\sqrt{wt}} \leq \frac{1}{(L+2\tau)(1+2\lambda)} \quad \Rightarrow \quad \lambda \leq \frac{1}{2} \left[ \frac{\sqrt{wt}}{L+2\tau} - 1 \right].$$

## C BACKGROUND ON MONOTONE OPERATORS

The results in our review can be found in different forms in standard textbooks (Rockafellar & Wets, 1998; Boyd & Vandenberghe, 2004; Nesterov, 2004; Bauschke & Combettes, 2017), and we include these results for completeness.

**Definition 1.** An operator  $\mathbf{T}$  is Lipschitz continuous with constant  $L > 0$  if

$$\|\mathbf{T}\mathbf{x} - \mathbf{T}\mathbf{y}\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n.$$

When  $L = 1$ , we say that  $\mathbf{T}$  is nonexpansive. When  $L < 1$ , we say that  $\mathbf{T}$  is a contraction.

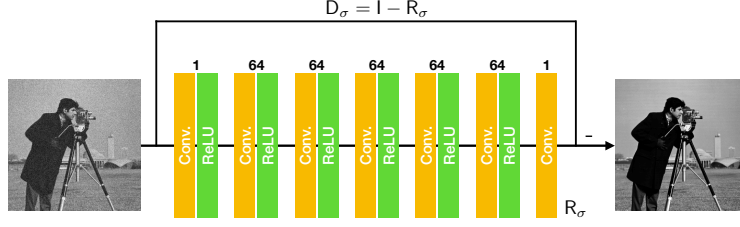


Figure 5: Illustration of the architecture of DnCNN used in all experiments. The neural net is trained to remove the AWGN from its noisy input image. We also constrain the Lipschitz constant of  $R_\sigma$  to be smaller than 2 by using the spectral normalization technique in (Sedghi et al., 2019). This provides a necessary condition for the satisfaction of Assumption 4.



Figure 6: Six test images used in the experiments on CS. From the left to right, there are *cameraman*, *house*, *pepper*, *starfish*, *butterfly*, and *jet*.

**Definition 2.**  $T$  is monotone if

$$(T(x) - T(y))^T(x - y) \geq 0, \quad x, y \in \mathbb{R}^n.$$

We say that it is strongly monotone or coercive with parameter  $\mu > 0$  if

$$(T(x) - T(y))^T(x - y) \geq \mu \|x - y\|^2, \quad x, y \in \mathbb{R}^n.$$

**Definition 3.**  $T$  is cocoercive with constant  $\beta > 0$  if

$$(T(x) - T(y))^T(x - y) \geq \beta \|Tx - Ty\|^2, \quad x, y \in \mathbb{R}^n.$$

When  $\beta = 1$ , we say that  $T$  is firmly nonexpansive.

The following results are derived from the definition above.

**Proposition 1.** For a convex and continuously differentiable function  $f$ , we have

$$\nabla f \text{ is } L\text{-Lipschitz continuous} \Leftrightarrow \nabla f \text{ is } (1/L)\text{-cocoercive}.$$

*Proof.* The proof is a minor variation of the one presented as Theorem 2.1.5 in Section 2.1 of (Nesterov, 2004).  $\square$

**Proposition 2.** Consider  $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$  and  $\beta > 0$ . Then, the following are equivalent

$$T \text{ is } \beta\text{-cocoercive} \Leftrightarrow I - 2\beta T \text{ is nonexpansive}.$$

*Proof.* Let  $R := I - 2\beta T$ , then  $T = 1/(2\beta)(I - R)$ . First suppose that  $T$  is  $\beta$ -cocoercive. Let  $h := x - y$  for any  $x, y \in \mathbb{R}^n$ . We then have

$$\beta \|T(x) - T(y)\|^2 \leq (T(x) - T(y))^T h = \frac{1}{2\beta} \|h\|^2 - \frac{1}{2\beta} (R(x) - R(y))^T h.$$

We also have that

$$\beta \|T(x) - T(y)\|^2 = \frac{1}{4\beta} \|h\|^2 - \frac{1}{2\beta} (R(x) - R(y))^T h + \frac{1}{4\beta} \|R(x) - R(y)\|^2.$$

By combining these two and simplifying the expression

$$\|R(x) - R(y)\| \leq \|h\|.$$

The converse can be proved by following this logic in reverse.  $\square$

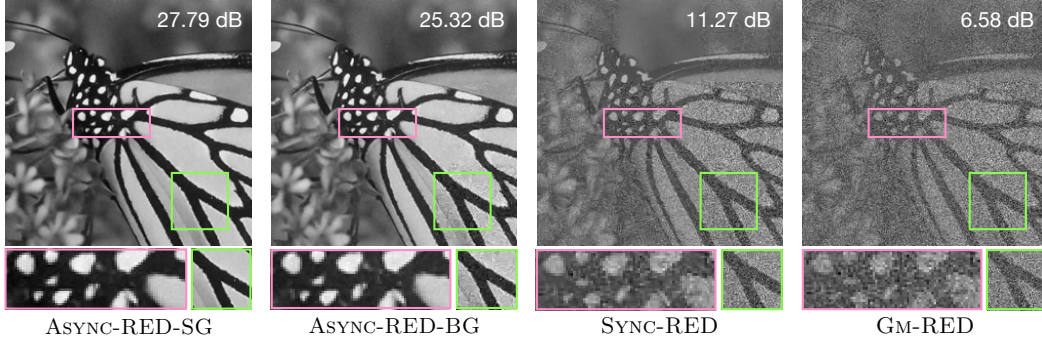


Figure 7: Visualization of the recovered images from the compressed measurements by ASYNC-RED-BG/SG, SYNC-RED, and GM-RED. Each algorithm is run with a time budget of 700 seconds.

The following characterization is also convenient.

**Proposition 3.** *For nonexpansive operators  $T_1$  and  $T_2$  with a constant  $\alpha \in (0, 1)$ , then the convex combination of the two operators  $(1 - \alpha)T_1 + \alpha T_2$  is nonexpansive.*

*Proof.* Let  $T := (1 - \alpha)T_1 + \alpha T_2$ . For any  $x, y \in \mathbb{R}^n$ , we can write

$$\|T(x) - T(y)\| \leq (1 - \alpha)\|T_1(x) - T_1(y)\| + \alpha\|T_2(x) - T_2(y)\| \leq \|x - y\|$$

□

## D ADDITIONAL TECHNICAL DETAILS

This section presents several technical details that were omitted from the main paper for space. Section D.1 presents the architecture and training of our DnCNN prior. Section D.2 provides extra details and validations that compliment the experiments in Section 5 of the main paper.

### D.1 ARCHITECTURE AND TRAINING OF THE DNCNN PRIOR

Our denoiser follows the standard architecture of DnCNN (Zhang et al., 2017a). Fig. 5 visualizes the architectural details of the DnCNN prior used in our experiments. Similar priors are extensively used in various PnP and RED algorithms (Zhang et al., 2017b; Ryu et al., 2019; Sun et al., 2019a). In total, the network contains 7 layers, of which the first 6 layers consist of a convolutional layer and a rectified linear unit (ReLU), while the last layer contains only a convolution operation. A skip connection from the input to the output is used to enforce the residual network  $R_\sigma$  to predict the noise residual. The output images of the first 6 layers have 64 feature maps, while that of the last layer is a single-channel image. We set all convolutional kernels to be  $3 \times 3$  with stride 1, which indicates that intermediate images have the same spatial size as the input image. We generated 44700 training examples by adding AWGN to 400 images from the BSD400 dataset (Martin et al., 2001) and extracting small patches of  $128 \times 128$  pixels with stride 30. Our DnCNN denoiser is trained to optimize the *mean squared error* by using the Adam optimizer (Kingma & Ba, 2015).

Different approaches have been used to constrain the Lipschitz constant (LC) of the denoising prior (Ryu et al., 2019; Sun et al., 2019a). We adopt the spectral normalization technique in (Sedghi et al., 2019) to control the LC of our DnCNN prior. In the training, we constrain the residual network  $R_\sigma$  such that its LC is smaller than 2. Since the non-expansiveness of  $D_\sigma$  implies that  $R_\sigma$  has  $LC \leq 2$ , this provides a *necessary* condition for  $D_\sigma$  to satisfy Assumption 4 (Sun et al., 2019a).

### D.2 EXTRA DETAILS AND VALIDATIONS

All experiments are run on the server equipped with 32 Intel(R) Xeon(R) CPU E5-2620 v4 processors of 3.2 GHz and 264 GBs of DDR memory. We trained all neural nets using NVIDIA RTX 2080

Table 1: SNR values obtained by ASYNC-RED-BG using different *block sizes* on CS task.

Block size	cameraman	house	pepper	starfish	butterfly	jet	Average
120	27.77	30.92	29.60	28.23	28.89	28.90	<b>29.05</b>
80	27.75	30.95	29.58	28.28	28.78	28.76	<b>29.01</b>
60	27.74	30.95	29.65	28.12	28.85	28.71	<b>28.99</b>

Table 2: SNR values obtained by ASYNC-RED-SG using different *minibatch sizes* on CS task.

minibatch size	cameraman	house	pepper	starfish	butterfly	jet	Average
1120	27.00	30.56	28.99	26.98	27.83	27.09	<b>28.03</b>
2240	27.31	30.78	29.01	27.64	28.29	28.03	<b>28.51</b>
3360	27.56	30.87	29.53	28.24	28.72	28.71	<b>28.93</b>

GPUs. We define the SNR (dB) used in the experiments as

$$\text{SNR}(\hat{\mathbf{x}}, \mathbf{x}) \triangleq 20 \log_{10} \left( \frac{\|\mathbf{x}\|_2}{\|\mathbf{x} - \hat{\mathbf{x}}\|_2} \right)$$

where  $\hat{\mathbf{x}}$  represents the reconstructed image and  $\mathbf{x}$  denotes the ground truth. Note that our experimental setup satisfies Assumption [1](#)[3](#) but provides a necessary condition for Assumption [4](#).

Fig. [6](#) shows the six test images used in the experiments of CS. They are resized to the size of  $240 \times 240$  pixels by using the Matlab function `imresize`. As demonstrated in the middle figure in Fig. [3](#), ASYNC-RED-SG converges faster than ASYNC-RED-BG given a fixed amount of time. This is further visualized in Fig. [7](#) where each algorithm is run for roughly 700 seconds. Since ASYNC-RED-SG uses only one-fourth of the total measurements, the per-iteration complexity is lower than ASYNC-RED-BG, leading to the faster convergence speed. In particular, the final SNR value obtained by ASYNC-RED-SG is roughly 2 dB higher than ASYNC-RED-BG. Additionally, both ASYNC-RED-BG/SG achieves significantly better results than SYNC-RED and GM-RED due to their adoption of asynchronous updates.

Table [1](#) and [2](#) illustrate the evolution of the reconstruction performance as the block size  $b$  and minibatch size  $w$  changes, respectively. Table [1](#) summarizes the SNR values obtained for three block sizes  $b \in \{60, 80, 120\}$ . Async-RED-BG achieves almost the same SNR values under these settings. Table [2](#) summarizes the SNR values for three minibatch sizes  $w \in \{1120, 2240, 3360\}$ , which corresponds to  $1/4$ ,  $1/2$ , and  $3/4$  of the full batch. As  $w$  increases, the final SNR performance improves, which is consistent with our theory. On the other hand, the error term due to stochastic processing in Theorem [2](#) is also proportional to the step size  $\gamma$ , which means that by using smaller  $\gamma$ , Async-RED-SG can approximate GM-RED as accurately as desired. However, a reduction in  $\gamma$  would also lead to slower convergence. One thus needs to tradeoff the desired accuracy against the desired speed to select a suitable configuration for Async-RED-SG.

The benefit of ASYNC-RED is fully explored when the denoiser acts like a *block-wise denoiser*, which means that it can perform denoising on blocks as effective as on the full image. A simple strategy for making denoisers block-effective is to include additional neighboring pixels at the input, but use the exact block size at the output. Table [3](#) reports the results of experimenting with the idea of *input padding* for DnCNN. The results indicate that by including a small number of pixels around each block at the input of DnCNN, one can match the performance of using the full image at the input of DnCNN.

The test image used in the experiment of CT is selected from the dataset of human protein atlas ([Williams et al., 2017](#)). We download 51 images that have the size of  $3000 \times 3000$  pixels. We select one image for test, which is cropped to  $800 \times 800$  pixels. We extract 39000 patches

Table 3: SNR values obtained by ASYNC-RED-BG using different *pad size* on CS task.

Pad size	cameraman	house	pepper	starfish	butterfly	jet	Average
0	27.64	30.86	29.44	28.09	28.64	28.69	<b>28.89</b>
10	27.72	30.99	29.56	28.25	28.79	28.73	<b>29.00</b>
20	27.72	30.99	29.57	28.27	28.79	28.75	<b>29.01</b>
full	27.75	30.95	29.58	28.28	28.78	28.76	<b>29.01</b>

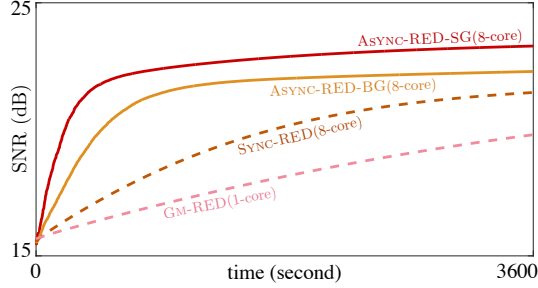


Figure 8: Convergence Illustration of ASYNC-RED-BG/SG and GM-RED for CT reconstruction with a time budget of 1 hour. Here, ASYNC-RED-SG randomly uses one-third of the total measurements at every iteration.

from the rest 50 images to train five specific DnCNN denoisers for the removal of AWGN with  $\sigma \in \{5, 10, 15, 20, 25\}$ . We report the result that has the highest SNR values. The Radon matrix used in the experiments corresponds to 180 angles with 1131 detectors. We synthesize the measurements by multiplying the Radon matrix with the vectorized image and add AWGN corresponding to 70 dB input SNR. In all tests, ASYNC-RED-SG randomly uses the measurements of 60 angles at each iteration, while ASYNC-RED-BG uses the entire measurement set. Fig. 8 plots SNR against the iteration number for Async-RED-BG/SG, Async-RED, and Gm-RED. Due to the lower per-iteration complexity, Async-RED-SG achieves the highest SNR value within the time budget of 1 hour. Fig. 9 provides a corresponding visual comparison between these methods. As reference, we also include the proximal gradient method with total variation regularizer (PGM-TV). The visual result of each method is obtained by running the algorithm with a time budget of 1 hour. Specifically, the per-iteration time cost of ASYNC-RED-BG/SG, SYNC-RED, GM-RED, PGM-TV are 5.23, 3.21, and 13.13, 19.19, and 44.74 seconds, respectively. The results clearly demonstrate that ASYNC-RED are indeed effective and efficient for a realistic, nontrivial imaging task on a large-scale image.

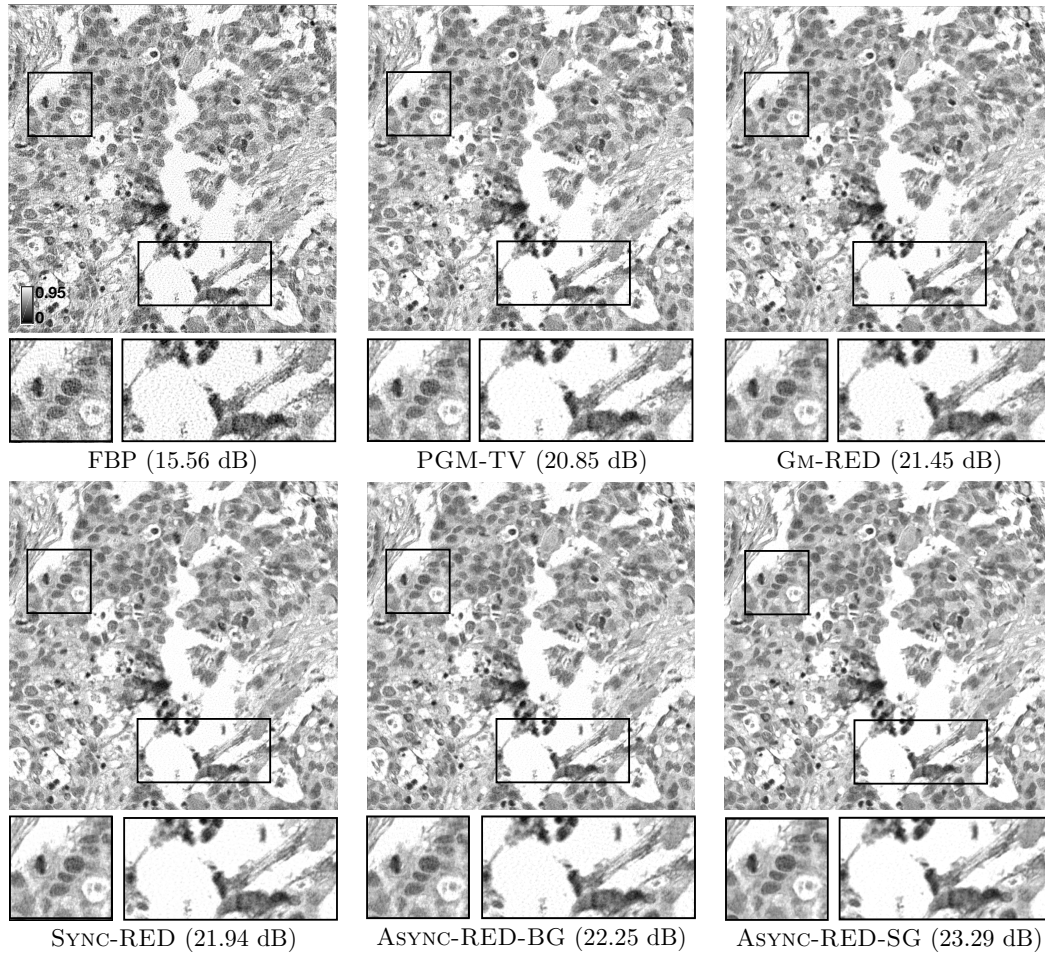


Figure 9: Visualization of the reconstructed CT images by PGM-TV, GM-RED, SYNC-RED, and ASYNC-RED-BG/SG. Each algorithm is run with a time budget of 1 hour. The colormap is adjusted for the best visual quality.