
Supplementary Materials

Anonymous Author(s)

Affiliation

Address

email

1 Missing Proofs

In this section, we will prove the main theorem in our paper. Before the detailed proof, we first recall the following assumptions that are commonly used for characterizing the convergence of non-convex stochastic optimization.

Assumption 1. (Bounded Gradient.) It exists $G \geq 0$ s.t. $\|\nabla f(\mathbf{w})\| \leq G$.

Assumption 2. (Bounded Variance.) It exists $\sigma \geq 0$ s.t. $\mathbb{E}[\|g(\mathbf{w}) - \nabla f(\mathbf{w})\|^2] \leq \sigma^2$.

Assumption 3. (L-smoothness.) It exists $L > 0$ s.t. $\|\nabla f(\mathbf{w}) - \nabla f(\mathbf{v})\| \leq L\|\mathbf{w} - \mathbf{v}\|, \forall \mathbf{w}, \mathbf{v} \in \mathbb{R}^d$.

1.1 Proof of Theorem 1

Based on the objective function of Sharpness-Aware Minimization (SAM), suppose we can obtain the noisy observation gradient $g(\mathbf{w})$ of true gradient $\nabla f(\mathbf{w})$, we can write the iteration of SAM:

$$\begin{cases} \mathbf{w}_{t+\frac{1}{2}} = \mathbf{w}_t + \rho \cdot \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|} \\ \mathbf{w}_{t+1} = \mathbf{w}_t - \eta \cdot g(\mathbf{w}_{t+\frac{1}{2}}) \end{cases} \quad (1)$$

Lemma 1. For any $\rho > 0, L > 0$ and the differentiable function f , we have the following inequality:

$$\langle \nabla f(\mathbf{w}_t), \nabla f(\mathbf{w}_t + \rho \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|}) \rangle \geq \|\nabla f(\mathbf{w}_t)\|^2 - \rho LG$$

Proof. We first add and subtract a term $\|\nabla f(\mathbf{w}_t)\|$ to make use of classical inequalities bounding $\langle \nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2), \mathbf{w}_1 - \mathbf{w}_2 \rangle$ by $\|\mathbf{w}_1 - \mathbf{w}_2\|^2$ for smooth.

$$\begin{aligned} LHS &= \langle \nabla f(\mathbf{w}_t), \nabla f(\mathbf{w}_t + \rho \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|}) - \nabla f(\mathbf{w}_t) \rangle + \|\nabla f(\mathbf{w}_t)\|^2 \\ &= \frac{\|\nabla f(\mathbf{w}_t)\|}{\rho} \langle \frac{\rho}{\|\nabla f(\mathbf{w}_t)\|} \nabla f(\mathbf{w}_t), \nabla f(\mathbf{w}_t + \rho \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|}) - \nabla f(\mathbf{w}_t) \rangle + \|\nabla f(\mathbf{w}_t)\|^2 \\ &\geq -L \frac{\|\nabla f(\mathbf{w}_t)\|}{\rho} \|\frac{\rho}{\|\nabla f(\mathbf{w}_t)\|} \nabla f(\mathbf{w}_t)\|^2 + \|\nabla f(\mathbf{w}_t)\|^2 \\ &= -\|\nabla f(\mathbf{w}_t)\| \rho L + \|\nabla f(\mathbf{w}_t)\|^2 \\ &\geq -G \rho L + \|\nabla f(\mathbf{w}_t)\|^2 \end{aligned}$$

where the first inequality is that

$$\langle \nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2), \mathbf{w}_1 - \mathbf{w}_2 \rangle \geq -L\|\mathbf{w}_1 - \mathbf{w}_2\|^2,$$

and the second inequality is the Assumption 1. \square

16 **Lemma 2.** For $\rho > 0$, $L > 0$, the iteration 1 satisfies following inequality:

$$\mathbb{E}\langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \geq \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 - L^2 \rho^2 - L\rho G$$

17 *Proof.* We denote the deterministic values of $\mathbf{w}_{t+\frac{1}{2}}$ as $\hat{\mathbf{w}}_{t+\frac{1}{2}} = \mathbf{w}_t + \rho \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|}$ in this section.

18 After we add and subtract the term $g(\hat{\mathbf{w}}_{t+\frac{1}{2}})$, we have the following equation:

$$\langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle = \langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_t + \rho \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|}) - g(\hat{\mathbf{w}}_{t+\frac{1}{2}}) \rangle + \langle g(\hat{\mathbf{w}}_{t+\frac{1}{2}}), \nabla f(\mathbf{w}_t) \rangle$$

19 For the first term, we bound it by using the smoothness of $g(\mathbf{w})$:

$$\begin{aligned} -\langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_t + \rho \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|}) - g(\hat{\mathbf{w}}_{t+\frac{1}{2}}) \rangle &\leq \frac{1}{2} \|g(\mathbf{w}_t + \rho \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|}) - g(\hat{\mathbf{w}}_{t+\frac{1}{2}})\|^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &\leq \frac{L^2}{2} \|\mathbf{w}_t + \rho \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|} - \hat{\mathbf{w}}_{t+\frac{1}{2}}\|^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &= \frac{L^2}{2} \|\mathbf{w}_t + \rho \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|} - (\mathbf{w}_t + \rho \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|})\|^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &= \frac{L^2 \rho^2}{2} \|\frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|} - \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|}\|^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &\leq L^2 \rho^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \end{aligned}$$

20 For the second term, by using the Lemma 1, we have:

$$\begin{aligned} \mathbb{E}\langle g(\hat{\mathbf{w}}_{t+\frac{1}{2}}), \nabla f(\mathbf{w}_t) \rangle &= \langle \nabla f(\hat{\mathbf{w}}_{t+\frac{1}{2}}), \nabla f(\mathbf{w}_t) \rangle \\ &= \langle \nabla f(\mathbf{w}_t + \rho \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|}), \nabla f(\mathbf{w}_t) \rangle \\ &\geq \|\nabla f(\mathbf{w}_t)\|^2 - \rho LG \end{aligned}$$

21 Assembling the two inequalities yields to the result. □

22 **Lemma 3.** For $\eta \leq \frac{1}{L}$, the iteration 1 satisfies for all $t > 0$:

$$\mathbb{E}f(\mathbf{w}_{t+1}) \leq \mathbb{E}f(\mathbf{w}_t) - \frac{\eta}{2} \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \sigma^2 + \eta L^2 \rho^2 + (1 - L\eta)\eta LG\rho$$

23 *Proof.* By the smoothness of the function f , we obtain

$$\begin{aligned} f(\mathbf{w}_{t+1}) &\leq f(\mathbf{w}_t) - \eta \langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle + \frac{L\eta^2}{2} \|g(\mathbf{w}_{t+\frac{1}{2}})\|^2 \\ &= f(\mathbf{w}_t) - \eta \langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle + \frac{L\eta^2}{2} (\|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_{t+\frac{1}{2}})\|^2 - \|\nabla f(\mathbf{w}_t)\|^2 + 2\langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle) \\ &= f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_{t+\frac{1}{2}})\|^2 - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \\ &\leq f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + L\eta^2 \|g(\mathbf{w}_t) - g(\mathbf{w}_{t+\frac{1}{2}})\|^2 \\ &\quad - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \\ &\leq f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + L\eta^2 L^2 \|\mathbf{w}_t - \mathbf{w}_{t+\frac{1}{2}}\|^2 \\ &\quad - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \\ &= f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + \eta^2 L^3 \rho^2 \\ &\quad - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \end{aligned}$$

24 Taking the expectation and using Lemma 2 we obtain

$$\begin{aligned}
\mathbb{E}f(\mathbf{w}_{t+1}) &\leq \mathbb{E}f(\mathbf{w}_t) - \frac{L\eta^2}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2\mathbb{E}\|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + \eta^2L^3\rho^2 \\
&\quad - (1 - L\eta)\eta\mathbb{E}\langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \\
&\leq \mathbb{E}f(\mathbf{w}_t) - \frac{L\eta^2}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2\sigma^2 + \eta^2L^3\rho^2 \\
&\quad - (1 - L\eta)\eta\mathbb{E}\langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \\
&\leq \mathbb{E}f(\mathbf{w}_t) - \frac{L\eta^2}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2\sigma^2 + \eta^2L^3\rho^2 \\
&\quad - (1 - L\eta)\eta \left[\frac{1}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 - L^2\rho^2 - L\rho G \right] \\
&= \mathbb{E}f(\mathbf{w}_t) - \frac{\eta}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2\sigma^2 + \eta^2L^3\rho^2 + (1 - L\eta)\eta L^2\rho^2 + (1 - L\eta)\eta L\rho G \\
&= \mathbb{E}f(\mathbf{w}_t) - \frac{\eta}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2\sigma^2 + \eta L^2\rho^2 + (1 - L\eta)\eta LG\rho
\end{aligned}$$

25

□

26 **Proposition 1.** Let $\eta_t = \frac{\eta_0}{\sqrt{t}}$ and perturbation amplitude ρ decay with square root of t , e.g., $\rho_t = \frac{\rho_0}{\sqrt{t}}$.
27 For $\rho_0 \leq G\eta_0$ and $\eta_0 \leq \frac{1}{L}$, we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 \leq C_1 \frac{1}{\sqrt{T}} + C_2 \frac{\log T}{\sqrt{T}},$$

28 where $C_1 = \frac{2}{\eta_0}(f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T))$ and $C_2 = 2(L\sigma^2\eta_0 + LG\rho_0)$.

29 *Proof.* By Lemma 3, we replace ρ and η with $\rho_t = \frac{\rho_0}{\sqrt{t}}$ and $\eta_t = \frac{\eta_0}{\sqrt{t}}$, we have

$$\mathbb{E}f(\mathbf{w}_{t+1}) \leq \mathbb{E}f(\mathbf{w}_t) - \frac{\eta_t}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta_t^2\sigma^2 + \eta_tL^2\rho_t^2 + (1 - L\eta_t)\eta_tLG\rho_t.$$

30 Take telescope sum, we have

$$\mathbb{E}f(\mathbf{w}_T) - f(\mathbf{w}_0) \leq - \sum_{t=1}^T \frac{\eta_t}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \sum_{t=1}^T \frac{1}{t} + (L^2\eta_0\rho_0^2 - L^2G\eta_0^2\rho_0) \sum_{t=1}^T \frac{1}{t^{\frac{3}{2}}}$$

31 Under $\rho_0 \leq G\eta_0$, the last term will be less than 0, which means:

$$\mathbb{E}f(\mathbf{w}_T) - f(\mathbf{w}_0) \leq - \sum_{t=1}^T \frac{\eta_t}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \sum_{t=1}^T \frac{1}{t}.$$

32 With

$$\frac{\eta_T}{2} \sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 \leq \sum_{t=1}^T \frac{\eta_t}{2}\mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 \leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \sum_{t=1}^T \frac{1}{t},$$

33 we have

$$\begin{aligned}
\frac{\eta_0}{2\sqrt{T}} \sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 &\leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \sum_{t=1}^T \frac{1}{t} \\
&\leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \log T.
\end{aligned}$$

34 Finally, we achieve the result:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 \leq \frac{2 \cdot (f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T))}{\eta_0} \frac{1}{\sqrt{T}} + 2(L\sigma^2\eta_0 + LG\rho_0) \frac{\log T}{\sqrt{T}},$$

35 which shows that SAM can converge at the rate of $O(\log T/\sqrt{T})$. □

36 **1.2 Proof of Theorem 2**

37 Suppose we can obtain the noisy observation gradient $g(\mathbf{w}_t)$ of true gradient $\nabla f(\mathbf{w}_t)$, and the mask
38 \mathbf{m} , we can write the iteration of SAM: Consider the iteration of Sparse SAM:

$$\begin{cases} \tilde{\mathbf{w}}_{t+\frac{1}{2}} = \mathbf{w}_t + \rho \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|} \odot \mathbf{m}_t \\ \mathbf{w}_{t+1} = \mathbf{w}_t - g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \end{cases} \quad (2)$$

39 Let us denote the difference as $\mathbf{w}_{t+\frac{1}{2}} - \tilde{\mathbf{w}}_{t+\frac{1}{2}} = \mathbf{e}_t$.

40 **Lemma 4.** *With $\rho > 0$, we have:*

$$\mathbb{E}\langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle \geq \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 - 2L^2\rho^2 - L\rho G - L^2\|\mathbf{e}_t\|^2$$

41 *Proof.* Similar to Lemma 2, We denote the true gradient as $\hat{\mathbf{w}}_{t+\frac{1}{2}} = \mathbf{w}_t + \rho \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|}$, and also add
42 and subtract the item $g(\tilde{\mathbf{w}}_{t+\frac{1}{2}})$:

$$\langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle = \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) - g(\hat{\mathbf{w}}_{t+\frac{1}{2}}) \rangle + \langle g(\hat{\mathbf{w}}_{t+\frac{1}{2}}), \nabla f(\mathbf{w}_t) \rangle$$

43 For the first term, we bound it by using the smoothness of $g(\mathbf{w})$:

$$\begin{aligned} -\langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) - g(\hat{\mathbf{w}}_{t+\frac{1}{2}}) \rangle &\leq \frac{1}{2} \|g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) - g(\hat{\mathbf{w}}_{t+\frac{1}{2}})\|^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &\leq \frac{L^2}{2} \|\tilde{\mathbf{w}}_{t+\frac{1}{2}} - \hat{\mathbf{w}}_{t+\frac{1}{2}}\|^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &= \frac{L^2}{2} \|\mathbf{w}_{t+\frac{1}{2}} - \mathbf{e}_t - \hat{\mathbf{w}}_{t+\frac{1}{2}}\|^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &\leq L^2 (\|\mathbf{w}_{t+\frac{1}{2}} - \hat{\mathbf{w}}_{t+\frac{1}{2}}\|^2 + \|\mathbf{e}_t\|^2) + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &= L^2 (\rho^2 \left\| \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|} - \frac{\nabla f(\mathbf{w}_t)}{\|\nabla f(\mathbf{w}_t)\|} \right\|^2 + \|\mathbf{e}_t\|^2) + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \\ &\leq 2L^2\rho^2 + L^2\|\mathbf{e}_t\|^2 + \frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 \end{aligned}$$

44 For the second term, we do the same in Lemma 2:

$$\mathbb{E}\langle g(\hat{\mathbf{w}}_{t+\frac{1}{2}}), \nabla f(\mathbf{w}_t) \rangle \geq \|\nabla f(\mathbf{w}_t)\|^2 - \rho LG.$$

45 Assembling the two inequalities yields to the result. □

46 **Lemma 5.** *For $\eta \leq \frac{1}{L}$, the iteration 2 satisfies for all $t > 0$:*

$$\begin{aligned} \mathbb{E}f(\mathbf{w}_{t+1}) &\leq \mathbb{E}f(\mathbf{w}_t) - \frac{\eta}{2} \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2\sigma^2 + 2\eta L^2\rho^2 + (1 - L\eta)\eta LG\rho \\ &\quad + (1 + L\eta)\eta L^2\|\mathbf{e}_t\|^2 \end{aligned}$$

47 *Proof.* By the smoothness of the function f , we obtain

$$\begin{aligned}
f(\mathbf{w}_{t+1}) &\leq f(\mathbf{w}_t) - \eta \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle + \frac{L\eta^2}{2} \|g(\tilde{\mathbf{w}}_{t+\frac{1}{2}})\|^2 \\
&= f(\mathbf{w}_t) - \eta \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle + \frac{L\eta^2}{2} (\|\nabla f(\mathbf{w}_t) - g(\tilde{\mathbf{w}}_{t+\frac{1}{2}})\|^2 - \|\nabla f(\mathbf{w}_t)\|^2 + 2\langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle) \\
&= f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t) - g(\tilde{\mathbf{w}}_{t+\frac{1}{2}})\|^2 - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle \\
&\leq f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + L\eta^2 \|g(\mathbf{w}_t) - g(\tilde{\mathbf{w}}_{t+\frac{1}{2}})\|^2 \\
&\quad - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle \\
&\leq f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + L\eta^2 L^2 \|\mathbf{w}_t - \tilde{\mathbf{w}}_{t+\frac{1}{2}}\|^2 \\
&\quad - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle \\
&= f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + \eta^2 L^3 \|\mathbf{w}_t - \mathbf{w}_{t+\frac{1}{2}} + \mathbf{e}_t\|^2 \\
&\quad - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle \\
&\leq f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + 2\eta^2 L^3 (\|\mathbf{w}_t - \mathbf{w}_{t+\frac{1}{2}}\|^2 + \|\mathbf{e}_t\|^2) \\
&\quad - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle \\
&= f(\mathbf{w}_t) - \frac{L\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + 2\eta^2 L^3 (\rho^2 + \|\mathbf{e}_t\|^2) \\
&\quad - (1 - L\eta)\eta \langle \nabla f(\mathbf{w}_t), g(\tilde{\mathbf{w}}_{t+\frac{1}{2}}) \rangle
\end{aligned}$$

48 Taking the expectation and using Lemma 4 we obtain

$$\begin{aligned}
\mathbb{E}f(\mathbf{w}_{t+1}) &\leq \mathbb{E}f(\mathbf{w}_t) - \frac{L\eta^2}{2} \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \mathbb{E}\|\nabla f(\mathbf{w}_t) - g(\mathbf{w}_t)\|^2 + 2\eta^2 L^3 (\rho^2 + \|\mathbf{e}_t\|^2) \\
&\quad - (1 - L\eta)\eta \mathbb{E}\langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \\
&\leq \mathbb{E}f(\mathbf{w}_t) - \frac{L\eta^2}{2} \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \sigma^2 + 2\eta^2 L^3 (\rho^2 + \|\mathbf{e}_t\|^2) \\
&\quad - (1 - L\eta)\eta \mathbb{E}\langle \nabla f(\mathbf{w}_t), g(\mathbf{w}_{t+\frac{1}{2}}) \rangle \\
&\leq \mathbb{E}f(\mathbf{w}_t) - \frac{L\eta^2}{2} \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \sigma^2 + 2\eta^2 L^3 (\rho^2 + \|\mathbf{e}_t\|^2) \\
&\quad - (1 - L\eta)\eta \left[\frac{1}{2} \|\nabla f(\mathbf{w}_t)\|^2 - 2L^2 \rho^2 - L\rho G - L^2 \|\mathbf{e}_t\|^2 \right] \\
&= \mathbb{E}f(\mathbf{w}_t) - \frac{\eta}{2} \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \sigma^2 + 2\eta L^2 \rho^2 + (1 - L\eta)\eta LG\rho \\
&\quad + (1 + L\eta)\eta L^2 \left\| \rho \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|} \odot \mathbf{m}_t - \rho \frac{g(\mathbf{w}_t)}{\|g(\mathbf{w}_t)\|} \right\|^2 \\
&= \mathbb{E}f(\mathbf{w}_t) - \frac{\eta}{2} \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta^2 \sigma^2 + 2\eta L^2 \rho^2 + (1 - L\eta)\eta LG\rho \\
&\quad + (1 + L\eta)\eta L^2 \|\mathbf{e}_t\|^2
\end{aligned}$$

49

□

50 **Proposition 2.** Let us $\eta_t = \frac{\eta_0}{\sqrt{t}}$ and perturbation amplitude ρ decay with square root of t , e.g.,
51 $\rho_t = \frac{\rho_0}{\sqrt{t}}$. With $\rho_0 \leq G\eta_0/2$, we have:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 \leq C_3 \frac{1}{\sqrt{T}} + C_4 \frac{\log T}{\sqrt{T}},$$

52 where $C_3 = \frac{2}{\eta_0} (f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + \eta_0 L^2 \rho^2 (1 + \eta_0 L) \frac{\pi^2}{6})$ and $C_4 = 2(L\sigma^2 \eta_0 + LG\rho_0)$.

53 *Proof.* By taking the expectation and using Lemma 4, and taking the schedule to be $\eta_t = \frac{\eta_0}{\sqrt{t}}$,
 54 $\rho_t = \frac{\rho_0}{\sqrt{t}}$, we obtain:

$$\begin{aligned} \mathbb{E}f(\mathbf{w}_{t+1}) &\leq \mathbb{E}f(\mathbf{w}_t) - \frac{\eta_t}{2} \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 + L\eta_t^2\sigma^2 + 2\eta_t L^2\rho_t^2 + (1 - L\eta_t)\eta_t LG\rho_t \\ &\quad + (1 + L\eta_t)\eta_t L^2\|\mathbf{e}_t\|^2 \end{aligned}$$

55 By taking sum and bound ρ with $\frac{G\eta_0}{2}$, we have:

$$\begin{aligned} \frac{\eta_0}{2\sqrt{T}} \sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 &\leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \sum_{t=1}^T \frac{1}{t} \\ &\quad + \sum_{t=1}^T (1 + L\eta_t)\eta_t L^2\|\mathbf{e}_t\|^2 \\ &\leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \sum_{t=1}^T \frac{1}{t} \\ &\quad + \eta_0 L^2\rho_0^2 \sum_{t=1}^T \frac{1}{t^{\frac{3}{2}}} + \eta_0^2 L^3\rho_0^2 \sum_{t=1}^T \frac{1}{t^2} \\ &\leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \sum_{t=1}^T \frac{1}{t} \\ &\quad + \eta_0 L^2\rho_0^2 \sum_{t=1}^T \frac{1}{t^2} + \eta_0^2 L^3\rho_0^2 \sum_{t=1}^T \frac{1}{t^2} \\ &\leq f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + (L\sigma^2\eta_0^2 + LG\rho_0\eta_0) \log T \\ &\quad + \eta_0 L^2\rho_0^2(1 + \eta_0 L) \frac{\pi^2}{6} \end{aligned}$$

56 Finally, we achieve the result:

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\nabla f(\mathbf{w}_t)\|^2 &\leq \frac{2(f(\mathbf{w}_0) - \mathbb{E}f(\mathbf{w}_T) + \eta_0 L^2\rho_0^2(1 + \eta_0 L) \frac{\pi^2}{6})}{\eta_0} \frac{1}{\sqrt{T}} \\ &\quad + 2(L\sigma^2\eta_0 + LG\rho_0) \frac{\log T}{\sqrt{T}} \end{aligned}$$

57 □

58 So far, we have completed the proof of the theory in the main text.

59 **2 More Experiments**

60 **VGG on CIFAR10.** To further confirm the model-agnostic characteristic of our Sparse SAM, we test
 61 the VGG-style architecture on CIFAR10. Following [?], we test SSAM training the VGG11-BN on
 62 CIFAR10 and the results are shown in the following Table 1. The perturbation magnitude ρ is set to
 63 0.05.

Table 1: Test accuracy of VGG11-BN on CIFAR10 with proposed Sparse SAM.

Model	Dataset	Optimizer	Sparsity	Accuracy	
VGG11-BN	CIFAR10	SGD	1	93.42%	
		SAM	0%	93.87%	
		SSAM-F/SSAM-D	50%	94.03%	93.79%
			80%	93.83%	93.95%
			90%	93.76%	93.85%
			95%	93.77%	93.48%
			98%	93.54%	93.54%
99%	93.47%	93.33%			

64 **SAM with different perturbation magnitude ρ .** We determine the perturbation magnitude ρ by
 65 using grid search. We choose ρ from the set $\{0.01, 0.02, 0.05, 0.1, 0.2, 0.5\}$ for CIFAR, and choose ρ
 66 from $\{0.01, 0.02, 0.05, 0.07, 0.1, 0.2\}$ for ImageNet. We show the results when varying ρ in Table 2
 67 and Table 3. From this table, we can see that the $\rho = 0.1$, $\rho = 0.2$ and $\rho = 0.07$ is suitable for
 68 CIFAR10, CIFAR100 and ImageNet respectively.

Table 2: Test accuracy of ResNet18 and WideResNet28-10 on CIFAR10 and CIFAR100 with different perturbation magnitude ρ .

Dataset	SAM ρ	0.01	0.02	0.05	0.1	0.2	0.5
CIFAR10	ResNet18	96.58%	96.54%	96.68%	96.83%	96.32%	93.16%
	WideResNet28-10	97.26%	97.34%	97.31%	97.48%	97.29%	95.13%
CIFAR100	ResNet18	79.56%	79.98%	80.71%	80.65%	81.03%	77.57%
	WideResNet28-10	82.25%	83.04%	83.47%	83.47%	84.20%	84.03%

Table 3: Test accuracy of ResNet50 on ImageNet with different perturbation magnitude ρ .

datasets	SAM ρ	0.01	0.02	0.05	0.07	0.1	0.2
ImageNet	ResNet50	76.63%	76.78%	77.12%	77.25%	77.00%	76.37%

69 **Ablations of Masking Strategy.** For further verification of our masking strategy, we perform more
 70 ablations in this paragraph. For the mask update in SSAM-F, the parameters with largest fisher
 71 information are selected. Compared with SSAM-F, we consider the random mask, *i.e.*, the mask is
 72 randomly generated to choose which parameters are perturbed. For the mask update in SSAM-D,
 73 we first drop the flattest weights and then random grow some weights. Compared with SSAM-D,
 74 we experiment the SSAM-D which drops randomly or drops the sharpest weights, *i.e.*, the weights
 75 with large gradients. The results of ablations are shown in Table 4. The results show that random
 76 strategies are less effective than our SSAM. The performance of SSAM-D dropping sharpest weights
 77 drops a lot even worse than random strategy, which is consistent with our conjecture.

78 **Influence of hyper-parameters** We first examine the effect of the number of sample size N_F
 79 of SSAM-F in Table 5. From it we can see that a certain number of samples is enough for the
 80 approximation of data distribution in SSAM-F, *e.g.*, $N_F = 128$, which greatly saves the computational
 81 cost of SSAM-F. In Table 6, we also report the influence of the mask update interval on SSAM-F and
 82 SSAM-D. The results show that the performance degrades as the interval becomes longer, suggesting
 83 that dense mask updates are necessary for our methods. Both of them are ResNet18 on CIFAR10.

Table 4: Ablation of different masking strategy.

Model	Dataset	Optimizer	Strategy	Accuracy
ResNet50	ImageNet	SGD	/	76.67%
		SAM	/	77.25%
		Sparse SAM	Random Mask	77.08%(-0.17)
		SSAM-F	Topk Fisher Information	77.31%(+0.06)
		SSAM-D	Random Drop	77.08%(-0.17)
			Drop Sharpnest weights	76.68%(-0.57)
	Drop flattest weights	77.25%(-0.00)		

Table 5: Results of ResNet18 on CIFAR10 with different number of samples N_F in SSAM-F. ‘Time’ reported in table is the time cost to calculate Fisher Information based on N_F samples.

Sparsity	N_F	Acc	Time
0.5	16	96.77%	1.49s
	128	96.84%	4.40s
	512	96.67%	15.35s
	1024	96.83%	30.99s
	2048	96.68%	56.23s
	4096	96.66%	109.31s
0.9	16	96.79%	1.47s
	128	96.50%	5.42s
	512	96.43%	15.57s
	1024	96.75%	29.24s
	2048	96.62%	57.72s
	4096	96.59%	110.65s

Table 6: Results of ResNet18 on CIFAR10 with different T_m intervals of update mask. The left of ‘/’ is accuracy of SSAM-F, while the right is SSAM-D.

Sparsity	T_m	Acc
0.5	1	96.81%/96.74%
	2	96.51%/96.74%
	5	96.83%/96.60%
	10	96.71%/96.73%
	50	96.65%/96.75%
	Fixed	96.57%/96.52%
0.9	1	96.70%/96.65%
	2	93.75%/96.63%
	5	96.51%/96.69%
	10	96.67%/96.74%
	50	96.64%/96.66%
	Fixed	96.21%/96.46%

84 3 Limitation and Societal Impacts

85 **Limitation.** Our method Sparse SAM is mainly based on sparse operation. At present, the sparse
86 operation that has been implemented is only 2:4 sparse operation. The 2:4 sparse operation requires
87 that there are at most two non-zero values in four contiguous memory, which does not hold for us. To
88 sum up, there is currently no concrete implemented sparse operation to achieve training acceleration.
89 But in the future, with the development of hardware for sparse operation, our method has great
90 potential to achieve truly training acceleration.

91 **Societal Impacts.** In this paper, we provide a Sparse SAM algorithm that reduces computation burden
92 and improves model generalization. In the future, we believe that with the development of deep
93 learning, more and more models need the guarantee of generalization and also the efficient training.
94 Different from the work on sparse networks, our proposed Sparse SAM does not compress the model
95 for hardware limited device, but instead accelerates model training. It’s helpful for individuals or
96 laboratories which are lack computing resources.