

1 We thank the reviewers for their time and helpful suggestions, which we will use to improve the paper’s presentation.
2 **R1, R2 (relevance of interpolation):** Please refer to lines 39-42 for examples where interpolation is satisfied. Recent
3 work [5,7,42,46,77] views interpolation as key to understanding the effectiveness of SGD for deep learning. Moreover,
4 we utilized this assumption to make algorithmic contributions that result in better empirical performance.
5 **R2 (comparison with [5], [77]):** They both use a constant step-size of $\frac{1}{L}$, which is either unknown or gives an
6 overly-conservative, small step-size. Our initial experiments confirmed that it lead to worse empirical performance and
7 we will mention this.
8 **R2, R3 (wall-clock):** For the line-search, we did ensure that the number of additional function evaluations is not large
9 (Section 7). In the Fig. 1 below, we show the wall-clock time per iteration averaged across training for the three datasets.
10 **R3 (error bars):** The figures in Section 7.3 do have error bars, but they unfortunately look like spikes in the submitted
11 version. We include one figure below with clearer error bars and will similarly update the remaining figures.
12 **R3 (Fixing typo for Theorem 1):** We correct the proof and statement of Theorem 1 below. Starting from the line
13 justified by Equation 2 in Appendix B (recall that $\mu_{ik} = 0$ if the f_{ik} is not strongly-convex),

$$\mathbb{E} \left[\|w_{k+1} - w^*\|^2 \right] \leq \left(1 - \mathbb{E}_{ik} \left[\mu_{ik} \min \left\{ \frac{1}{L_{ik}}, \eta_{\max} \right\} \right] \right) \|w_k - w^*\|^2$$

14 We consider the following two cases: $\eta_{\max} < 1/L_{\max}$ and $\eta_{\max} \geq 1/L_{\max}$. When $\eta_{\max} < 1/L_{\max}$,

$$\mathbb{E} \left[\|w_{k+1} - w^*\|^2 \right] \leq (1 - \mathbb{E}_{ik} [\mu_{ik} \eta_{\max}]) \|w_k - w^*\|^2 = (1 - \bar{\mu} \eta_{\max}) \|w_k - w^*\|^2$$

15 When $\eta_{\max} \geq 1/L_{\max}$, we use $\min \left\{ \frac{1}{L_{ik}}, \eta_{\max} \right\} \geq \min \left\{ \frac{1}{L_{\max}}, \eta_{\max} \right\}$ to obtain

$$\mathbb{E} \left[\|w_{k+1} - w^*\|^2 \right] \leq \left(1 - \mathbb{E}_{ik} \left[\mu_{ik} \frac{1}{L_{\max}} \right] \right) \|w_k - w^*\|^2 = \left(1 - \frac{\bar{\mu}}{L_{\max}} \right) \|w_k - w^*\|^2.$$

16 Combining the two cases gives us the theorem statement with L_{max} instead of L . We will make this change in
17 Theorem 1 statement. Note that Theorem 4’s proof will be changed similarly.

18 **R3 (Requested) rigorous proof for Theorem 3:** We can prove an $O(1/T)$ rate by bounding $\eta_{\max} \leq \frac{3}{2\rho L}$ as follows:

$$\frac{f(w_{k+1}) - f(w_k)}{\eta_k} \leq \frac{L\eta_k}{2} \|\nabla f_{ik}(w_k)\|^2 - \langle \nabla f(w_k), \nabla f_{ik}(w_k) \rangle \quad (\text{Using smoothness and dividing by } \eta_k)$$

$$\implies \mathbb{E} \left[\frac{f(w_{k+1}) - f(w_k)}{\eta_k} \right] \leq \left(\frac{L\eta_{\max}\rho}{2} - 1 \right) \|\nabla f(w_k)\|^2 \quad (\text{Since } \eta_k \leq \eta_{\max} \text{ and using the SGC})$$

$$\|\nabla f(w_k)\|^2 \leq \frac{1}{1 - \frac{L\eta_{\max}\rho}{2}} \mathbb{E} \left[\frac{f(w_k) - f(w_{k+1})}{\eta_k} \right] \quad (\text{Rearranging and upper-bounding } \eta_{\max} \leq \frac{2}{L\rho}.)$$

$$\implies \|\nabla f(w_k)\|^2 \leq \left(\frac{1}{1 - \frac{L\eta_{\max}\rho}{2}} \right) \left(\frac{1}{\eta_{\max}} + \frac{L_{max}}{2(1-c)} \right) \mathbb{E} [f(w_k) - f(w_{k+1})].$$

(Bounding η_k using the line-search similar to Appendix C)

19 Telescoping and setting $c = 1/2$ and $\eta_{\max} \leq \frac{3}{2\rho L}$ completes the proof. It is non-trivial to avoid the dependence of ρ, L in
 η_{\max} and we leave it as future work. Regardless of this result, we believe that this paper’s contributions are impactful.

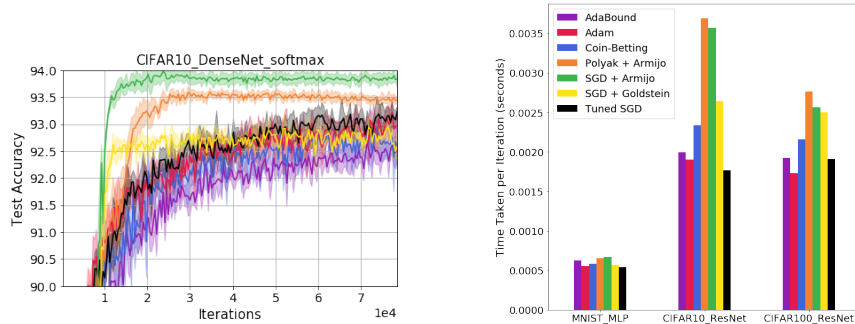


Figure 1: **Left:** CIFAR-10 with new error-bar style. **Right:** Average iteration times on CIFAR-10.