

1 We would like to thank the four reviewers for their feedback. We first discuss the common concern about our
2 contributions and novelty shared by [reviewer 2](#), [reviewer 3](#) as follows.

3 **•• Key Contributions & Novelty:** While the algorithm design idea of combining SPIDER and EM can be natural, this
4 paper involves a significant work in the convergence analysis technique and gives new insights to deriving future
5 stochastic EM algorithms. We highlight four particular challenges and contributions: (A) **EM is not a gradient-based**
6 **algorithm**, while SPIDER is designed for accelerating gradient algorithm. The convergence analysis of EM-based
7 algorithms can not be deduced from that of gradient-based procedures: they require a specific study. As a key observation
8 to outline this difference, note that EM can be equivalently studied in the parameter space and in the expectation space -
9 which is not the case for the gradient methods. (B) Among the incremental-EM techniques, we provide **state-of-the-art**
10 complexity bounds that overpass all the previous ones. (C) Contrary to previous incremental EM algorithms [6,15],
11 the approximate mean field H_{k+1} is **biased** except at the beginning of the epoch (Lemma 4, Corollary 5). A main
12 novelty of the proof is therefore to show that the squared error when approximating the mean field, is of the same
13 order as the squared of the approximation (Proposition 6) - contrary to previous EM-based proofs, the usual martingale
14 increment property no longer holds, thus requiring a novel approach for the control of this error. (D) We provide a **new**
15 **perspective** to interpret SPIDER-EM by showing its equivalence to a perturbed onLine-EM where the perturbation acts
16 as a control variate in order to reduce the variance (see Section 5).

17 **Reviewer 1.** We thank the reviewer for the careful review of our paper. Please find our response as follows.

18 **Linear Convergence:** The reviewer has made a mindful observation. Indeed, in Fig. 3, SPIDER-EM appears to converge
19 linearly, i.e., faster than our theory’s prediction. We speculate that this is due to a local strong convexity condition of
20 $W(\cdot)$ around a local minimum statistics \hat{s} , after SPIDER-EM reaches a neighborhood of \hat{s} . This condition was never
21 made in our analysis, however it has been used as an assumption in sEM-vr *without verification*. Nevertheless, under
22 this condition, we can show local linear convergence for SPIDER-EM. The proof idea will be included in the final
23 version. Lastly, Balakrishnan et al. consider a non-incremental, first order EM that optimizes θ via gradient ascent, this
24 is different from SPIDER-EM that operates on the statistics (\hat{S}), this reference will be included in the final version.

25 **Complexity Guarantee with $b = \mathcal{O}(1)$:** This can be easily derived from Theorem 2. If $b \asymp n^a$, $k_{in} \asymp n^c$ in the case
26 $c \geq a \geq 0$, then the overall complexity will be $K_{CE}(n, \epsilon) = n + \mathcal{O}(n^{\max\{\frac{a+c}{2}, 1-\frac{a+c}{2}\}} \epsilon^{-1})$, $K_{Opt}(n, \epsilon) = \mathcal{O}(n^{\frac{c-a}{2}} \epsilon^{-1})$.
27 Setting $b = \mathcal{O}(1)$ or equivalently $a = 0$ will lead to a tradeoff in complexity in terms of the number of parameter
28 updates and of conditional expectations computed. We will include a discussion for these general settings.

29 **Comparison to Karimi et al.:** The quasi-gradient interpretation of the E-step updates (see Proposition 1 and H5) is
30 similar to those from Karimi et al. which analyzed the onLine EM, and the same technique has also been used in [8,15].
31 Our analysis are different as there are additional complexity due to the *biased, variance reduced* estimator, note that
32 this improves the convergence rate from $\mathcal{O}(1/\epsilon^2)$ (onLine EM) to $\mathcal{O}(\sqrt{n}/\epsilon)$ (SPIDER-EM). That said, Karimi et al. is
33 relevant to this work, which we will discuss in the final version.

34 **Second Order Stationary (SOS) point:** This is an interesting question. The difficulty in applying the analysis for
35 SOS point from [10] lies on the non-gradient update nature of SPIDER-EM. In fact, the procedures in Sec. 3.3. from
36 [10] hinge on the availability of approximate Hessian. To our best knowledge, the convergence to SOS point **has not**
37 **been studied** in the context of incremental EM. A possible solution is to apply the Louis’s missing info. principle [8]
38 to estimate the Hessian of the log-likelihood, and thus diagnosing if the stationary point is SOS. We will include a
39 discussion about the challenge in the final version.

40 **Reviewer 2:** Thank you for the careful review of our paper. As per your suggestion, we will include a table in the final
41 paper that compares the complexity results to existing works [also see the response on **Complexity with $b = \mathcal{O}(1)$**].

42 We emphasize that our contributions are beyond *an application of SPIDER on EM algorithm*. Instead, as explained in
43 the beginning of this rebuttal, our analysis involve non-trivial proof techniques to bound the error of approximation to
44 the mean field (non-gradient) update of SPIDER-EM. This analysis is different from sEM-vr, FIEM in [6,15] (inspired by
45 SVRG, SAGA) while the approximate mean field for SPIDER-EM is **biased**. With such challenges, it is not obvious that
46 the analysis of SPIDER can be directly applied to SPIDER-EM. In fact, our algorithm design/analysis is only **inspired**
47 by SPIDER, for it involves a brand new analysis on incremental-EM algorithms.

48 **Reviewer 3:** Thank you for the careful review and positive feedback. We have addressed your comments regarding the
49 contributions and novelty in the beginning of this rebuttal. In the final version, we will include a highlight about them.

50 **Reviewer 4:** Thank you for the careful review and positive feedback. Here are the specific answers to the comments.

51 **Assumptions H4.3 and H5:** These assumptions are actually classical in the analysis of incremental/stochastic EM
52 algorithms - a relevant reference is [15], where explicit examples satisfying H4.3, H5 are provided. In fact, $B(s)$ can be
53 written in terms of the Hessian of the function in H3, which is positive definite if the map $T(s)$ is unique. As shown in
54 [15], this holds for a number of distributions such as GMM when $R(\theta)$ is a properly designed barrier function.

55 **Parameter/Statistics Space:** We emphasize that it is essential for the EM algorithm to work alternatively in the
56 expectation space (E-step) and the parameter space (M-step). However, we agree that the present writing is not optimal
57 and we will try to streamline the discussions in the final version.