1 We would like to thank three reviewers for their feedback. Upon acceptance, we will include in the final version (a)
2 *new local linear convergence results for fiEM method*, (b) *an improved presentation of main results* and (c) *missing*
3 *references*. We first discuss a few common concerns shared by **reviewer 1**, **reviewer 2**, **reviewer 3**.

4 • • **Local Linear Convergence of fiEM**: As observed by the reviewers, empirically fiEM shows a local linear
5 convergence similar to sEM-VR. We found that fiEM has local linear convergence **in theory**, too. The new analysis
6 requires same assumptions as [Thm.1,Chen+2018] and adopts proof of [Defazio+2014]. We show $\mathbb{E}\|\hat{\mathbf{s}}^{(k)} - \mathbf{s}^\star\|^2 \leq$
7 $(1 - \delta)^k \|\hat{\mathbf{s}}^{(0)} - \mathbf{s}^\star\|^2$ for $k \geq 0$ with $\delta = \Theta(1/n)$, where $\mathbf{s}^\star$ is a stationary point to (19).

8 • • **Satisfaction of Assumptions**: <u>All</u> assumptions H1-H5 are verified rigorously in the GMM, pLSA applications
9 presented, as proven in Appendix G. They are mild even though should be checked on a case-by-case basis. Reviewers
10 are referred to [McLachlan&Krishnan 2007] which shows satisfaction of similar assumptions on a variety of applications.

11 • • **Clarity:** We admit it is a challenging task to present all technical results within the page limit, but we will try our
12 best to improve in the final version, viz. using a running example to illustrate the assumptions used and implementation
13 of algorithms. We will also clarify about the expectation operators in theorems and correct typos.

14 **Reviewer 1:** We thank the reviewer for valuable comments and references. Our point-to-point response is as follows:

15 **Related work:** The paper [Karimi+2019] is relevant and will be included. Thank you for bringing it to our attention.
16 Karimi+[2019] focused on a biased stochastic approximation scheme and gave a global convergence rate for sEM. In
17 this case, their analysis shares similar scaled gradient interpretation as fiEM and sEM-VR, yet w/o variance reduction.

18 **iEM's Rate**: You are right as the rate of iEM is comparable to GD. Yet, iEM is a popular method without a previously
19 known global rate. Indeed, the comparison of iEM to fiEM, sEM-VR (theoretical & empirical) is our main contribution.

20 **Comparison to [Chen+2018]**: Our assumptions are more practical and less restrictive. Global convergence to stationary
21 point for sEM-VR in [Thm.2, Chen+2018] assumes i) the sufficient statistics $\mathbf{s}_i(\overline{\boldsymbol{\theta}}(\boldsymbol{s}'))$ is $L_s$-Lipschitz continuous in $\boldsymbol{s}'$,
22 $\forall i$ – this is implied by our H1-H5 via Lemma 4; ii) the complete log-likelihood is strongly concave – this is slightly
23 relaxed in our H4 which only requires a unique global minimizer for the complete log-likelihood. Besides, H1-H5 are
24 **directly verifiable** (as explained above) and we provide the rate towards a stationary point. Lastly, local convergence in
25 [Thm.1, Chen+2018] requires $\|\hat{\mathbf{s}}^{(k)} - \mathbf{s}^\star\|$ to be in a ball of radius $\mathcal{O}(1/L_s)$ for **any** $k \geq 1$. This is a strong assumption
26 that is not directly verifiable even if $\hat{\mathbf{s}}^{(0)} \approx \mathbf{s}^\star$ is known a-priori.

27 **Reviewer 2:** We thank the reviewer for useful comments. Please find the comparison of fiEM, sEM-VR below:

28 **Comparing fiEM to sEM-VR:** This comparison is analogous to comparing SAGA to SVRG for finite sum optimization,
29 and there is no clear ordering. In short, it depends on the trade-off of memory imprint and computation complexity.
30 sEM-VR requires $\mathcal{O}(\dim(\mathsf{S}))$ space to store $\overline{\mathbf{s}}^{(\ell(k))}$, yet a *full pass* on the data set is needed at each epoch, resulting in
31 higher complexity; meanwhile, fiEM only processes the data set *incrementally*, but it requires $\mathcal{O}(n\dim(\mathsf{S}))$ to store the
32 variables involved. We remark that the global rate for sEM-VR is **not proven** in [Chen+2018]. In Fig. 2 we show fiEM
33 outperformed sEM-VR in one dataset (a bigger one) but was outperformed by sEM-VR in the other (a smaller one).

34 **Reviewer 3:** We thank the reviewer for the comments. We clarify that *in addition* to analyzing iEM using the MISO
35 framework (which will be mentioned explicitly), we analyzed fiEM, sEM-VR with a **completely different** framework
36 w/ **scaled gradient**, the latter constitutes our main contribution of fast global convergence rates; see p.2 of our paper.

37 **Global Convergence & H4:** We emphasize H4 **does not** imply that every stationary point of (1) is global minima,
38 as having a unique global minimizer *does not* imply *any first order critical point is global minima*. Also, H4 refers
39 to *complete log-likelihood* $L(\boldsymbol{s}, \boldsymbol{\theta})$ with fixed $\boldsymbol{s}$, instead of incomplete log-likelihood $\overline{\mathcal{L}}(\boldsymbol{\theta})$ in (1). It holds for most
40 exponential family models where EM is useful [McLachlan&Krishnan 2007]. Mind that $\overline{\mathcal{L}}(\boldsymbol{\theta})$ is non-convex and our
41 convergence is *global* in the sense that it does not restrict the initialization, a common assumption for analysis of EM.

42 **Bounds in theorems:** The current presentation style of theorems, which evaluates the gradient norm of a ran-
43 domly terminated stochastic EM solution, is common in **stochastic non-convex** optimization e.g., [Ghadimi&Lan
44 2013,Reddi+2016a/b]. Part of the reason is that it results in a practical solution. While picking the best iterate leads to
45 the same sublinear rate as ours, doing so involves a full pass on the data ($\nabla\overline{\mathcal{L}}$) and computing the incomplete likelihood,
46 both are **difficult** tasks avoided in stochastic EM methods. Besides, as the reviewer mentioned, both random termination
47 and best iterate schemes lead to a quantity upper bounded by $\sum_k \mathbb{E}\|\nabla\overline{\mathcal{L}}(\boldsymbol{\theta}^{(k)})\|^2/K_{\max}$. This quantity is not equal to
48 the *averaged iterate*, and upper bounding it by $\mathcal{O}(1/K_{\max})$ is a non-trivial task – and is precisely our main contribution.

49 **Theorem 1 of Paper 1613:** Indeed, Theorem 1 for iEM is a special case of [Thm. 1, 1613]. We will cite the latter
50 properly. Our main contribution here lies on *fast convergence* of fiEM, sEM-VR shown by a different framework, see
51 Theorem 2. Detailed comments about the difference between this paper and 1613 has been sent to the AC.