

1 We thank the reviewers for their time, effort, and helpful feedback. We will make sure to correct all typos and incorporate  
2 the reviewers’ suggestions in the final version of the paper. We address individual feedback below.

3 **Reviewer 1:**

4 Yes, “stochastic noise” is the error in the gradient approximations, i.e.  $\nabla F(x^k) - g^k$  in the notation of the paper. We  
5 will add this definition to the paper.

6 Bounded noise is indeed a strong assumption, and we are not sure whether the proof can be extended to the more  
7 general case where  $\mathbf{E}_k[\|\xi^k\|^2]$  is uniformly bounded. We should note that [25] proves a similar result for SGM with a  
8 more general noise condition than bounded noise, and their technique may extend to QHM, but bounded noise greatly  
9 simplifies the proof. We will comment on this in the paper. In our opinion, the result showing convergence of QHM  
10 when  $\nu_k\beta_k \rightarrow 1$  is interesting from both theoretical and practical perspectives. From the theoretical side, the result  
11 shows that it is possible to always be increasing the amount of momentum (in the limit when  $\nu_k\beta_k = 1$ , we are not  
12 using the fresh gradient information at all!) and still obtain convergence (which is not too surprising, since the objective  
13 function is smooth). From the practical point of view, Theorem 5 shows that increasing  $\nu_k\beta_k \rightarrow 1$  might lead to smaller  
14 stationary distribution size, which may give better empirical results. We will include this discussion in the final version  
15 of the paper.

16 Theorem 3 guarantees that there exists some sequence  $\epsilon_k$  such that the statement holds. We will update the statement of  
17 Theorem 3 to include  $\mu > 0$  and  $\exists\{\epsilon_k\}$  with  $\epsilon_k \geq 0$ . When  $\beta = 1$  and is constant, the algorithm cannot converge, since  
18 the gradient information is not used at all. However, as we show in Section 2, it is possible to set  $\beta_k \rightarrow 1$  in the limit  
19 and still obtain convergence. As mentioned above, we found this difference theoretically interesting and will highlight  
20 it after Theorem 3.

21 We had to limit the explanation of Figures 1, 2, and 3 due to space constraints. We would use the additional page  
22 allowed in the final version to include more interpretation of these figures, explaining their relationships to Theorems 3,  
23 4, and 5 in more detail.

24 To improve the presentation of our theoretical results, we will include more lead-in discussion before our theorems  
25 (especially Theorem 3) to make them easier to understand. We believe that combining our literature review into its own  
26 section (as suggested by Reviewer 2 and addressed below) will also help the reader better contextualize our theorems.  
27 Regarding sensitivity, we will expand Section 5 with comments that the presented results depend continuously on the  
28 algorithm parameters.

29 **Reviewer 2:**

30 Since we presented a range of results covering different aspects of momentum methods (asymptotic convergence,  
31 stability, and stationary distributions), and each of these areas has its own rich literature, we put a brief review of the  
32 relevant works at the beginning of each section. To improve readability, we will combine these references in a separate  
33 section to give a narrative that ties the different areas together.

34 We cite Mandt et al. [16] for some empirical justification of our Assumption A3 (in the beginning of Section 3) and for  
35 proving an equivalent result to our Theorem 4 in the case of stochastic heavy ball (in the beginning of Section 4). We  
36 can make the existing citation after Theorem 4 more prominent.

37 By saying that “we use quadratic functions for ease of analysis” we did not mean to imply that it is straightforward to  
38 prove our results for a more general class of functions, although the basic principles remain similar. We will change the  
39 wording of this phrase in the final version of the paper to make it more specific. The results in Theorems 4 and 5 may  
40 hold approximately near a local minimum for nonconvex but sufficiently smooth functions, when the learning rate is  
41 sufficiently small. In that case, the function is intuitively almost quadratic locally, since the higher-order terms in the  
42 Taylor expansion around the optimum can be neglected.

43 For most of our experiments, the code is very minimal. To make it easier to reproduce our results, we will release the  
44 code for all of the experiments performed in this paper.

45 **Reviewer 3:**

46 We agree with the reviewer that a more general analysis is necessary to better understand the properties of momentum  
47 methods in the wild, especially for the nonconvex problems faced in deep learning practice. While our results are based  
48 on some restrictive assumptions, we hope they still provide valuable practical insights, and we believe they could serve  
49 as an important foundation for any future work extending the unified analysis of momentum methods to more general  
50 classes of functions.