

1 We would like to thank four reviewers for their feedback.

2 Upon acceptance, we will include in the final version (a) *improved notations* and (b) *an improved presentation of*
3 *related work*. We first discuss a few common concerns shared by **reviewer 1**, **reviewer 2**, **reviewer 3** and **reviewer 4**.

4 **•• Non convex bound:** As pointed out by several reviewers, the non convex bound does not clearly show a dependence
5 on the gradient prediction process. While being clear that in the convex case, predicting well the next gradient
6 theoretically improves the bound (see Eq (2)), the non convex bound at least set a comparable state-of-the-art rate (as
7 adam-type methods) in $\mathcal{O}(\sqrt{1/T})$. We acknowledge the need for a more precise consideration of the constants of both
8 our method and AMSGrad to highlight not only an empirical edge of the optimistic update but also a theoretical one.

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

10 **Convex regret bound:** For analysis purposes we presented the algorithm without projection step by assuming the
11 compact assumption **H1**. Of course, this assumption needs to be verified and we partially did it for a model of interest
12 that is a deep neural network, see Section 4.3. Adding projection steps is a neat idea to avoid having those issues but is
13 not common in non convex optimization analyses, see references [5, 9, 14, 38].

14 **Numerical example:** We thank the reviewer for their remark on the numerical runs. The main motivation behind those
15 plots is to show that adding an optimistic update to the vanilla AMSGrad actually speed up the convergence in terms of
16 both losses and accuracies. Given the well-known advantages of Adam-type methods as ADAM or AMSGrad, we did
17 not compare to slower methods "that does not have any of the extra features of AMSGrad" as written by the reviewer.

18 **Reviewer 2:** We thank the reviewer for valuable comments. A proofreading of the paper is being done as suggested
19 and we give the following clarification:

20 **Wall clock times comparison:** We agree with the reviewer with the heavy computation that our gradient prediction
21 process can represent. In the shown runs, and as precised Section 5.3, only $r = 5$ gradients are being used for the
22 extrapolation step. Both memory and wall clock time are lightly impacted. **To complete. Can we have the wall clock**
23 **times plots?**

24 **Reviewer 3:** We thank the reviewer for the thorough analysis. Our remarks are listed below:

25 **Gradient prediction algorithm:** We agree with the reviewer that a study of how well the gradient is predicted using
26 the current method would be impactful. The scope of our paper being the stochastic optimization method itself, we
27 invoked a simple but effective gradient prediction algorithm on the basis of reference [31] which shows great theoretical
28 and empirical acceleration using such extrapolation. Of course, there is room for improvement regarding that prediction
29 process and can be the object of further research papers.

30 **To check if [31] gives theoretical guarantees on how well the prediction is**

31 **Numerical evaluation:** It has been rightly noted that in Figure 3, the curves are still rising and thus convergence is not
32 attained yet. Though, for illustrative purposes, the main idea is to show how faster our method is in the first epochs.
33 The purpose of this method is not to achieve better generalization (*i.e.* reach better accuracies at convergence) but rather
34 to show how less epochs are needed to achieve similar results as baselines. The learning rates have been tuned over a
35 grid search and the best results have been reported. The choice of a constant learning rate was made to stick to our
36 theoretical results. Runs with exponential decay or step decay can also be done for completeness.

37 **Reviewer 4:** We thank the reviewer for valuable comments and typos. Our response is as follows:

38 **Numerical Experiments:** We only reported the average of the 5 runs but as the reviewer suggested we will report error
39 bars in the rebuttal version of the paper.

40 We agree with the fact that our method empirically show on Figure 2 and 3 better training performances (both in terms
41 of loss and accuracy) but we must note how comparable and most of the time better than the baselines our method
42 behaves at testing phase.