

1 We would like to thank all four reviewers for the close look given at the paper. We believe that the paper gained clarity
2 and readability by taking their comments into account. We now answer the different issues and comments.

3 **Introduction (REVIEWER 1 AND 2)** *Accelerated optimization*: this choice of wording was made to emphasize
4 the fact that optimizing a smoothed version of the non-smooth objective function would lead faster to a (better) result.
5 However, following your advice and to avoid confusion with the sense commonly given in the field, we will replace this
6 term by *faster optimization*.

7 *Problem formulation*: we will introduce the problem tackled more formally in the introduction, introducing the objective
8 function f and describing Table 1 as the error bounds on the norm of the difference between the objective function and
9 its resulted smoothed form.

10 *Poor generalization capabilities*: our sentence is misleading, "resulting in poor generalization capabilities [8]" will be
11 deleted.

12 **Widening of local minima (REVIEWER 4)** The widening of local minima was indeed not sufficiently developed,
13 and we have thus decided to add the following result to Sec. 2.3 (using the extra page of the camera-ready version).

14 **Proposition 1** (Widening of local minima). *Let $\mu > 0$ and $y \in \mathbb{R}^d$. If f is L -Lipschitz, then there exists $z \in \mathbb{R}^d$ such
15 that $\|y - z\| \leq \mu$ and, $\forall x \in \mathbb{R}^d$,*

$$f_{\gamma, \alpha}(x) \leq f(y) + \frac{L}{2\mu} \|x - z\|^2 + \frac{d}{2\alpha} \left(1 + \ln \left(1 + \frac{\alpha\mu L}{d}\right)\right), \quad (1)$$

16 where $\gamma = \min \left\{ \mu, \sqrt{\frac{\mu}{L\alpha}} \right\}$.

17 In other words, for every local minima $y \in \mathbb{R}^d$ of the objective function f , $f_{\gamma, \alpha}(x)$ will be small (i.e., approximately
18 $f(y)$) in a neighborhood of y of size μ (see Fig. 1.c for an example). As a consequence, a good but thin local minimum
19 will have its basin of attraction increased, and thus be easier to reach by GD even when the starting point x_0 is far from
20 the local minimum. For example, consider $f(x) = \min\{1, |x|\}$. Its gradient is zero for any $x \notin [-1, 1]$, which means
21 that GD initialized outside this region will be stationary. Moreover, for large smoothing parameters γ , RS will tend
22 to flatten the objective and thus lead to the same behavior. However, BMR (with a sufficiently large α) will create an
23 almost quadratic function in a region $x \in [-\mu, \mu]$, thus allowing GD to converge even when initialized at distance μ
24 from the origin.

25 **Experiments (REVIEWER 2, 3 AND 4)** *Adding the comparison against ME would have been helpful*: Moreau
26 Envelope is a great tool when the inner optimization problem (i.e., the proximal operator) can be solved efficiently,
27 ideally in closed form. Here we focus on settings where there is no closed form, thus directly using Moreau Envelope
28 would add a full optimization problem at each step. The results would eventually need to be analyzed w.r.t. the precision
29 achieved when solving the ME problem, adding a layer of complexity to the analysis. We chose not to follow that
30 direction as a first go, in order to illustrate the core concepts unequivocally.

31 *Optimization for DL, 10 dimensions is a bit disappointing*: Our goal was here to emphasize the improvements over RS
32 on simple functions. Scaling the method to big ML problems such as the training of deep neural networks and policy
33 optimization in Reinforcement Learning is currently investigating and left for future work.

34 *BMR presents some oscillations on Figure 3*: this is due to the exploratory nature of the method in addition to the
35 particular choice of α in this experiment, leading to an opportunistic behavior. Lowering α would decrease the
36 oscillations observed toward the end of the process.

37 **Typos and references (REVIEWER 1, 2 AND 3)** Thank you for pointing to the relevant works [1, 3]. We will
38 include them in the paper. We will also fix the reference [2] in the proofs. Finally, the bias in $O(\sqrt{K})$ is indeed a typo,
39 it will be modified to $O(1/\sqrt{K})$.

40 **General comments (REVIEWER 3)** Our method bridges the gap between two different smoothing methods com-
41 monly used in optimization, both having their advantages and drawbacks. It must not be viewed as a trivial combination
42 of ME and RS but as an interpolation between them. This interpolation allows to take the best of both worlds and leads
43 to faster optimization, as shown in the experiments.

44 References

45 [1] A. Beck and M. Teboulle. Smoothing and first order methods: A unified framework. *SIAM Journal on Optimization*, 2012.

46 [2] C. Forbes, M. Evans, N. Hastings, and B. Peacock. *Statistical distributions*. John Wiley & Sons, 2011.

47 [3] Y. Nesterov. Smooth minimization of non-smooth functions. *Mathematical programming*, 2005.