

1 We thank Reviewers (R) 1, 2 and 3 (who gave us marks 6, 7, and 8, respectively) for their pertinent remarks.

2 **R1+R2+R3. More details on monotone operators.** We will provide more details on the monotone operator frame-
3 work (for instance more elaboration around Eq. 6) allowing us to prove the theorems (we shall use the 9th page of
4 the camera ready if our paper is accepted). We already provided some intuitions in the appendix, e.g. in Appendix
5 C. Moreover, we want to recall that the definition of the space E is provided in Appendix C. We will move it to the
6 beginning of Section 4.1 for better understanding.

7 **R1+R3. Straightforward combination of existing techniques?** As R3 says, our work closes the theoretical gap
8 by providing the first (and so far the only) optimal first order algorithm for smooth strongly convex decentralized
9 optimization. We obtained this algorithm by “combining existing approaches”, but the combination was far from
10 straightforward.

11 • (Scaman et al.’17) obtained MSDA by simply applying Nesterov acceleration to the dual problem. We instead
12 build upon recent results on the minimization of strongly convex functions under linear constraints (by a first order
13 algorithm, without projecting on the constraints space), see (Salim et. al.’20). Surprisingly, there are only a few
14 algorithms that can solve such problems at a linear rate, and they were proposed only recently.

15 • APAPC is obtained by applying Nesterov acceleration to the generalized forward-backward algorithm (5) for a sum
16 of operators A and B (see page 5, line 208). Although we managed to do this, this was not an easy task (to say the
17 least), because Nesterov acceleration does not apply to general monotone operators. Even if it did, a naive approach
18 would lead to a sublinear rate $\mathcal{O}\left(\frac{1}{k^2}\right)$ because $A + B$ is not strongly monotone. On the contrary, we obtained an
19 accelerated linear rate (complexity $\mathcal{O}\left((\sqrt{\chi\kappa} + \chi) \log \frac{1}{\epsilon}\right)$) in Theorem 2, which requires careful and deep theoretical
20 analysis of APAPC. Finally, we had to carefully design our generalized Forward Backward algorithm by choosing the
21 space E (see Appendix C), its inner product, and the matrix P as functions of the gossip matrix W .

22 • In Appendix F we provided an algorithm provably optimal in “# of communication rounds”, without using Chebyshev
23 acceleration. The development and analysis of this method required substantial innovation, as we explain in the paper.

24 Finally, we believe that the *apparent* simplicity of our approach is due to us spending a lot of time making sure the
25 explanations are as intuitive as possible. Many of these intuitions only became clear to us after we have done the
26 analysis; and we provide them for the benefit of the reader. Hence, we view the simplicity as a strength!

27 **R1. Experiments.** The networks chosen for evaluating the decentralized method are the ones that were used in
28 (Scaman et al.’17): 10×10 grid and Erdős-Rényi random network with parameter $p = 0.06$. We have now added this
29 detail to the paper. **Regarding the wall-clock-time comparison:** The design of our experiments was very similar to
30 those in (Scaman et al.’17), who assumed that local gradient computation takes one unit of time and communication
31 with neighbors takes time τ . It’s easy to observe that in this case [wall clock time] = [# of gradient calls] + $\tau \times$
32 [# of communication rounds]. Scaman et al. (2017) used 2 regimes in their experiments: $\tau \gg 1$ (high communication
33 time) and $\tau \ll 1$ (low communication time), which more or less correspond to our plots “# of communication rounds”
34 and “# of gradient calls”. This controlled setup is sufficient to verify numerically that our theory (which expresses
35 the optimality of OPAPC in terms of “# of communication rounds” and “# of gradient calls”) has predictive power for
36 experiments. Note that our “# of communication rounds” and “# of gradient calls” plots give understanding of how
37 the algorithms will behave with any possible τ (even if we do not know it in practice). We will also produce plots
38 providing “wall clock time”, but these will be implementation dependent. Our focus here was not to produce highly
39 performing and fine-tuned software to be benchmarked in this way.

40 **R3. Minor comments.** We will put Table 1 in the main paper using the 9th page if the paper is accepted.

41 **R3. Open remark on optimality w.r.t. W .** In the line of works on optimal distributed algorithms by Scaman,
42 Hendrikx, Xiao, Bubeck, Bach and Massoulié, lower bounds are obtained by proving the existence of a “bad” gossip
43 matrix and “bad” functions that cannot be optimized faster than the lower bounds by any decentralized algorithm. This
44 includes the decentralized algorithms using the gossip matrix W . Therefore, if one pick a family of functions and a
45 gossip matrix W , they could be “bad” in the above sense, and one cannot beat the lower bounds by a decentralized
46 algorithm using W . However, we agree that, perhaps, the lower bounds theory for these distributed algorithms could
47 be a bit improved by providing lower bounds involving intrinsic properties of the graph, as in the centralized case
48 (there are indeed many gossip matrices for one graph).

49 **R3. Open remark w.r.t. real-life application.** Since OPAPC is practical and optimal, the use of OPAPC at scale is
50 definitely the next step in the study of OPAPC. Obviously, this goes beyond the scope of this paper, which contains
51 algorithm development, analysis and testing, but does not and was not supposed to have a software/system development
52 element. But obviously we are also very curious about its performance at scale, and plan to work on this in the future.
53 We are optimistic and confident that OPAPC will outperform existing approaches on average, as indicated by our
54 experiments.