1 We appreciate the valuable comments from the reviewers. We will revise them accordingly.

2 **Novelty of DAPG and Brief History of Decentralized Proximal Gradient Descent**

3 The novelty of DAPG is the main concern of reviewers because DAPG is closely related to Mudag [25]. However,
4 reviewers don't know and we don't emphasize that it is hard to extend decentralized gradient descent methods with
5 gradient tracking to proximal counterpart with the same linear convergence rates. Please note that, researchers take *five*
6 *years* to propose DPA [1], which is the first linearly convergent decentralized proximal gradient algorithm, after the
7 publication of EXTRA [19], the first decentralized gradient descent method of linear convergence rate.

8 PG-EXTRA [17], the follow-up work of EXTRA, extends EXTRA to the composite setting which has extra non-smooth
9 term, but only obtains the sub-linear convergence rate even for strongly convex $f_i(x)$. In the following years, different
10 decentralized gradient descent algorithms with gradient tracking were proposed such as [13], [14]. However, no
11 evidence shows that these algorithms can be extended to the composite setting but keeping linear convergence. Even the
12 recent work NIDS [8] can only achieve a sub-linear convergence rate for the composite setting. In fact, until five years
13 after EXTRA, DPA [1] is the first linearly convergent proximal gradient algorithm for decentralized optimization which
14 is published in Neurips 2019. Just as the title of [13]-'Harnessing smoothness to accelerate distributed optimization'
15 mentioned, gradient tracking tries to harness the smoothness of the function. However, due to non-smooth term, the
16 convergence analysis of gradient tracking based methods become much harder.

17 From above history, we can observe that, in the decentralized setting, extending the results of decentralized gradient
18 descent methods to their proximal counterparts is not an easy work. Thus, whether Mudag [25] can be extended to the
19 decentralized proximal setting is not obvious. In fact, just as pointed by Reviewer 2, to deal with the non-smoothness,
20 DAPG takes more consensus steps compared to Mudag. The proof of Lemma 2 is also totally different from the one of
21 Lemma 9 of Mudag. Note that, the proof of DAPG does not rely on the technique of DPA, either.

22 Furthermore, the results obtained by DAPG are total new for decentralized proximal algorithms. To the best of
23 our knowledge, DAPG is the first accelerated decentralized gradient descent which theoretically outperforms current
24 decentralized proximal algorithms.

25 Thus, the novelty and contribution of DAPG are substantial.

26 **Reviewer_1** : Thank you for the comments on figures, we will revise accordingly.

27 **Reviewer_2**

28 Q1:Incremental contribution to Mudag

29 A1: Just as mentioned at the beginning, it is not easy to extend decentralized gradient descent with gradient tracking
30 to the composite setting. Because of the non-smoothness of composite setting, DAPG has more consensus steps than
31 Mudag and the proof of Lemma 2 is totally different from the one of Lemma 9 of Mudag.

32 **Reviewer_3**

33 Q1: Experiment and Relation to prior work and Reference to APM-C

34 A1: We should set $K$ according to the condition number of graph. We will give more experiments on graphs of large
35 condition numbers in our revised paper. We will give more detailed comparison with Mudag and will add the reference
36 of APM-C in our revised paper.

37 **Reviewer_4**

38 Q1: Comparison with Accelerated version of decentralized proximal gradient descent methods?

39 A1: As mentioned at the beginning, DPA [1] is the first decentralized proximal gradient descent method with linear
40 convergence rate. NIDS is shown to achieve linear convergence rate in [24] which is the best convergence rate of
41 decentralized proximal gradient methods can achieve before our work. Work [14] is a decentralized accelerated gradient
42 descent method which can not deal with the non-smooth regularization term and no evidence shows that it can be
43 directly extended to deal with the non-smooth regularization term. To the best of our knowledge, DAPG is the first
44 decentralized accelerated proximal gradient descent method. Thus, we have compared DAPG with state-of-the-art
45 algorithms theoretically and empirically.

46 Q2: What will happen if all the FastMix steps in Alg.1 are replaced by standard (K rounds of) communication? It is
47 unclear what will happen $f_i(x)$ is $\nu$ strongly convex?

48 A2: By standard communication, the communication complexity of algorithm will depend on $1/(1 - \lambda_2(W))$ instead
49 of $1/\sqrt{1 - \lambda_2(W)}$. DAPG does not require $f_i(x)$ to be $\nu$ strongly convex. But if $f_i(x)$ is $\nu$ strongly convex, the results
50 of Theorem 1 of course still hold.