1   We thank each of the reviewers for their comments and suggestions.

2   **Clarification on sparsity regime**   (raised by $R_1$ and $R_2$).

3   In our submission, we used the upper-bound notation $a, b = O(\log n)$ which includes $O(1)$.

4   We did so because behaviour exhibited in different regimes is fundamentally dichotomous. We wanted to be succinct
5   in gathering all of the results in different regimes.

6   As seen from Line 63 for *large changes*, reliable testing is possible in the $O(1)$ sparsity regime - the comments
7   following Thm. 2 explore this. Nevertheless, for *small changes* (see Line 67, and Thm. 1), it is impossible to reliably
8   test with $a, b = O(1)$ - rather, $a, b = \Omega(\log n)$ is both necessary and sufficient for testing small changes. Overall, this
9   is one of the main messages of the paper: the SNR requirements for testing large and small changes are qualitatively
10   different.

11   Thus, to present all of our testing results, we must technically permit $a, b$ to vary from constant to logarithmic, even
12   though we are primarily focused on what is possible when $a, b = O(1)$. That said, we will revise the surrounding text
13   to provide a bit more context for this technical note.

14   In addition, our thanks to $R_2$ for bringing to attention our oversight in not citing Mossel, Neeman, and Sly here, which
15   we will amend.

16   **Practicality**   (raised by $R_2$ and $R_3$). First, we note that our approach can be naturally extended to more practical
17   settings (e.g., many communities, degree correction). For instance, for multi-community setting, given the number of
18   communities, one can again recover weakly (e.g., by leveraging SDP methods of Fei and Chen), and then adapt our
19   statistics in a straightforward manner. However, this will affect the SNR thresholds for testing, and fully characterising
20   the dependence on the number of communities, etc. is messy and non-trivial, enough to merit further work (mirroring
21   the development of the recovery literature).

22   Similarly, while we present preliminary exploration of this in the experiments section, validation of these methods and
23   models on real-world networks is a non-trivial task. We think that extensive pursuit of these here would distract from
24   the primary theoretical considerations of the paper.

25   We thank you each for bringing up these important questions, and we will include a discussion of these as directions
26   for future research.

27   **Constants**   (raised by $R_2$). Note that, as stated on line 157, each constant in the paper can be *explicitly bounded*,
28   although these bounds may not be the tightest possible. For instance, in the limit as $n \nearrow \infty$, the TST result for large
29   changes shows impossibility of reliable testing if $\Lambda < 1$, and also that $\Lambda > 4$ is sufficient for reliable testing. The
30   lower bound is explicitly discussed in the proofs, see Line 866 of the supplement. The upper bound follows from the
31   proof of Thm. 2 and the work of Mossel, Neeman, and Sly. Non-asymptotic results are also discussed.

32   From our perspective such gaps within a constant factor of each other is acceptable. Establishing exact constants in
33   the SBM can be quite challenging, and, even for recovery, these are only known in very particular settings, such as
34   exact recovery or weak recovery with distortion $n(1/2 - \varepsilon)$. This is certainly an interesting problem, one that likely
35   requires a dedicated effort to resolve. We will add a discussion of this question to a future work section.

36   **Balance**   (raised by $R_2$). We will try to highlight the balance assumption in the introduction, although we note the
37   use of 'balance' in the title precisely for this purpose. (As an aside, since submission we have extended the theorems
38   to unbalanced but linearly-sized communities, perhaps ameliorating this concern.)

39   **Lines 183,184**   (raised by $R_1$). This was meant to be a brief comment explaining why we switch, based on the
40   relative sizes of $a$ and $b$, from a test based on counts of edges across communities $N_a^{x_0}(G)$ to one based on counts of
41   edges within communities $N_w^{x_0}(G)$. Both of these counts are of roughly the same signal strength, but depending on
42   the regime ($a > b$ vs. $a \leq b$), one count will have a lower noise level. We will clarify this comment so that it reads
43   more smoothly.

44   **Typographical errors, and presentation suggestions.**   We thank the reviewers for pointing these out, especially $R_1$
45   for catching one in a definition (the sup should be over $d = 0$ *or* $d \geq s$)! We will correct these.