1 We thank all the reviewers for their time and effort.

2 [APPROPRIATENESS FOR NEURIPS] As pointed out by Reviewer 1, reviewing practices are of interest to the NeurIPS
3 community. Moreover, *our work is applicable beyond peer review as well, in applications such as admission or hiring*,
4 where the evaluation assignments are not random and where one wishes to avoid excessive evaluations for the testing.
5 We note that importantly, *past NeurIPS editions feature a number of works on statistical testing* (e.g., "Hypothesis
6 Testing in Unsupervised Domain Adaptation with Applications in Alzheimer's Disease" (2016) and "Differentially
7 Private Uniformly Most Powerful Tests for Binomial Data" (2018) to just name a few due to space constraints here),
8 and we believe our work on statistical hypothesis testing aligns well with the scope of NeurIPS.

9 [IMPORTANCE OF THE PROBLEM] It is extremely important to note that the claim "This review format is furthermore
10 not widely adopted anyway" does not hold throughout most of academia. Even in computer science, most Theory
11 conferences use single blind. There are massive ongoing debates on SB vs. DB in many fields of academia (including
12 in fields such as databases which have switched to DB, but there is some push to move back to SB). Based on our
13 experience in this debate, those supporting SB also argue for evidence of biases *in their specific community*. Indeed, an
14 important roadblock in making the SB vs. DB arguments is the absence of rigorous evidence – and this is where our
15 work significantly contributes by designing principled procedures to answer the pressing questions in these debates.

16 [EFFECT SIZES] We thank Reviewer 1 for raising an important point of effect sizes. In fact, we should have mentioned
17 that **the test statistic of our test represents the effect size**. This is in line with the seminal work by J. Cohen ("A
18 power primer", 1992.) where the test statistic is suggested to estimate the effect size for the sign test. We are more than
19 happy to add more in-depth discussion of the effect size in the revision.

20 [GLMS AND OTHER MODELS] Our test makes significantly fewer assumptions than the tests based on GLMs and other
21 parametric models. Specifically, our test *does not rely on strong modelling assumptions* (in contrast to GLMs) and *also
22 holds when reviewer decisions can be completely subjective* (in contrast to models that assume a presence of "true"
23 [latent] qualities). Indeed, the models suggested by the reviewers are at a risk of making spurious conclusions because
24 of the restrictive natures of the suggested models which may not capture the highly complex decision-making process
25 of the human reviewers.

26 [ASSIGNMENT AND MATCHING] During the review period, we **have solved the open problem of designing the
27 experimental procedure** that leads to a reliable testing. Specifically, we have designed the experimental procedure
28 (allocation of reviewers to conditions, assignment and matching) that follows standard conference peer-review pipeline
29 (i.e., allows to perform any assignment algorithm, including TPMS) and *does not inflate Type-I error rates of the testing
30 procedures*. Our proposed procedure assigns reviewers to conditions and papers to reviewers jointly in a carefully
31 selected manner, thereby avoiding issues pertinent to the setup of Tomkins et al. All the statements in the paper hold for
32 this procedure, but do not require an idealized random assignment any more. If reviewers prefer, we can include it in
33 the final version.

34 [PREVIOUS WORKS] Conference peer review setup is not a fully randomized controlled trial (i.e., the reviewers are
35 not assigned at random) and hence *past approaches fail due to idiosyncrasies of the process*. With respect to the
36 specific work mentioned by Reviewer 1 (Bertrand & Mullainathan, 2004), their method assigns identities of authors to
37 (fabricated) documents at random. In our setup, *random assignment of author identities to real (i.e., non-fabricated)
38 submissions* is problematic due to various logistical and ethical issues including reviewers guessing actual authors
39 thereby causing biases, not all authors/researchers agreeing to have their paper/name modified, and others — this
40 opens a separate can of worms which should be rigorously addressed before using it in the peer review setting. We are
41 definitely happy to add a discussion of this and other relevant papers (including those mentioned by Reviewer 3) in the
42 revision.

43 [PERMUTATION TEST] The standard way of performing the permutation test would fail to control for the Type-I
44 error because of the additional confounding due to quality of submissions. Our test is a careful modification of the
45 permutation test which provably controls for the Type-I error rate even in presence of such confoundings. In the final
46 version, we are happy to detail the shortcomings of the standard permutation test in our setup.

47 [REAL DATA EXPERIMENT] Unfortunately, the data from the Tomkins et al. experiment is not available to us. Tomkins
48 et al. mention in their work that releasing this data (even in an anonymized format) would make it possible to
49 deanonymize reviewers. Through our developed toolkit, we are happy to assist any program chairs who are interested to
50 conduct tests of biases in their respective research community.