

1 Thank you all for the thoughtful reviews! We will respond in plural form to maintain blindness. We will first address  
2 issues brought up by multiple reviewers, followed by individual queries.

3 **Regarding the definition of reproducibility:** We regarded a paper as reproducible if most (75%+) of the claims  
4 in the paper could be reproduced. If a claimed improvement was measured in orders-of-magnitude, being within  
5 the same order-of-magnitude was sufficient (e.g., a paper claims 700x faster, but we observe 300x, still qualifies).  
6 When compared to other algorithms, we consider a paper reproduced if most (90%+) of the new algorithm's rankings  
7 correspond to what was in the paper (e.g., new method was most accurate on 95% of tasks compared to 4 other models,  
8 we want to see our reproduction be most accurate on at least  $95% * 90% = 81%$  of the same tasks, compared to the same  
9 models). As a last resort, we considered getting within 10% of the numbers reporting in the paper (or better), or in the  
10 case of non-quantitative results (e.g., GAN sample quality), we subjectively compare our results with the paper to make  
11 a decision. We will include a version of the above in the revised paper.

12 **Regarding Venue:** While a "science of science" type venue would be appropriate, we feel NeurIPS is a more appropriate  
13 venue. We are more concerned with the nature of reproducibility within our specific field, rather than broader genres  
14 such as computer science or science generically. We feel the discussion appropriate given the current high and growing  
15 discussion of these issues within the field, and necessary to build communal momentum around larger organized efforts  
16 to track and quantify reproduction (as we discussed regarding the ICLR Reproducibility Challenge). While we consider  
17 our work valuable, issues in study bias will persist until we get a larger pool of implementers under consideration.

18 **Regarding Author/reproducer details:** We are trying to maintain double-blindness as much as possible in our replies.  
19 The camera ready will detail this and reproducer background extensively.

20 **To R#1:** We would have also liked to do a textual analysis of papers, and hope to do so in the future. However, first  
21 attempts produced significant issues with respect to parsing reliability (some papers are photo-copies, some are OCR'd,  
22 some are poorly formatted PDFs, some are nicely formatted PDFs, etc) that would confound results and increase analysis  
23 difficulty, not to mention take considerable time to get working reliably. This also prevented us from automating the  
24 equation counting.

25 Re line 212: you make an excellent point, and we will add caution to the final manuscript regarding that statement.

26 "*Furthermore, the reader would also benefit from a rough classification of reasons for non-reproducible papers*": We  
27 did not record this in detail for each paper, but agree would be valuable! We will add a discussion of the general types  
28 of issues we encountered in the final version (though unfortunately non-quantified). Subjectively, the below list would  
29 be the primary issues that gave us reproduction problems, which we will elaborate further in camera ready. None of the  
30 below issues were mutually exclusive.

- 31 1. Unclear notation or language. A component of the algorithm is explained, but not in a way easily understood  
32 by the reproducers, or was ambiguously specified.
- 33 2. Missing algorithm step or details. A step was completely left out of description.
- 34 3. Results left as an exercise to the reader: many papers would specify loss functions or other equations for which  
35 the gradient needed to be taken, but then not detail the resulting gradients. Depending on the functions and  
36 math involved re-deriving was non-trivial.
- 37 4. Missing hyper-parameters, or similar nuance details. We appear to have an implementation accurate to what  
38 was described, but some minor detail was not specified and makes a big difference in results.

39 We understand knowing the number of reproducers / backgrounds is intrinsically valuable for our paper, but also violates  
40 double blind review. Reproduction attempts/effort was approximately uniformly distributed between reproducers.

41 **To R#2:** The effort was indeed significant, and only possible because we used software early on that recorded much  
42 relevant information. Between efforts done as part of education, job, and for fun, back-of-envelope guesstimation puts  
43 our effort at  $\approx 10,710$  hours per author.

44 **To R#3:** On computing "*how long it took to reproduce*", this is a very good idea, thank you! We think we could do this  
45 for *most* of our papers that were reproduced, as most were made open source and we can compare our start date to the  
46 commit date to get an approximation. We can't do it for all though as some remain closed source, and we would be  
47 unable to compare to unreproduced papers since we did not track that information.

48 We can't get this done by the end of rebuttal time, but will consider it for camera ready - we aren't sure how long this  
49 would take us (given day jobs as well!). We also want to give it the same thought and consideration we gave the factors  
50 already in the paper, and not rush analysis without considering confounding factors.