

1 We thank all reviewers for their careful reviews and many positive comments, including **R1**: "really valuable to the
2 neuro community", "gives a roadmap for using NNs...to tell us how brains work", "I really liked this paper"; **R4**:
3 "this paper is novel and significant," "well-written and relatively easy to understand"; **R3**: "could be an interesting
4 contribution to the field." Even the most negative reviewer **R2** stated: "a highly original contribution with huge potential
5 in the field," "the analysis...is of uniformly high quality," "with the supplement, the methods are clear and most model
6 explanations make sense." We now address major reviewer concerns and clarify our contributions, as detailed below:

7 **Extension to deeper CNNs without spatial invariances in stimuli (R2,R3,R4)**: While we demonstrated a novel ap-
8 plication of attribution methods to model reduction of 1-hidden layer CNNs, specifically to validate deep CNNs in
9 the retina where we had neurobiological ground truth, we can easily extend our method to deeper CNNs of depth D
10 processing natural movies through a dynamic programming (DP) approach that works backwards from layer D to
11 layer 1. First, note a natural movie of limited duration without spatial invariances is still well approximated by a low
12 dimensional trajectory in both pixel space and every hidden layer. Let K be the max dimensionality for spatial input
13 patterns for any channel in any layer. Then the basic idea is to attribute the response in layer D to the K dimensional
14 space of inputs to each channel in layer $D - 1$ using integrated gradients. We first find the important channels in
15 layer $D - 1$ using methods in our paper. Then we recursively iterate via the same method to layer $D - 2$ and so
16 on down to the pixel layer. Because of the DP-like nature of our algorithm, the computational complexity (after
17 dimensionality reduction to K) is $O(DKC)$ where C is the max number of channels in a layer, and not exponential in
18 D as **R3** worried. The end result is a set of important channels in each layer, along with, for each important channel
19 $\leq K$ linear combinations of neurons that matter for generating the response in layer D . We are actively pursuing this
20 method in deeper networks, but we will share pseudocode for this algorithm in a revised version before acceptance to
21 NeurIPS. However, consistent with **R2,R3,R4** we feel completing this program is well beyond the scope of this paper,
22 especially since neurobiological ground truth is missing for higher areas. But we hope our success in the retina and the
23 extendability of our approach to deeper networks, will provide a great roadmap for neuroscience as recognized by **R1**.

24 **Experimental evidence for our new model of omitted stimulus response (OSR) (R3)**: As shown quantitatively in
25 [2], the model subunits match bipolar cells (BCs), and the 3 in the OSR correspond to fast OFF, fast ON and slow
26 ON BCs, thus mapping directly to biological pathways. Furthermore, multiple BC types can connect to a ganglion
27 cell (GC) (Asari and Meister 2012). Thus our new model is basically consistent with known anatomy. However, we
28 leave *further* physiological validation of our model, beyond successfully generating the OSR, to future work, which
29 would require painstaking experiments to perturb BC pathways and observe GC responses. We believe it is already a
30 substantial contribution to show our approach *automatically* extracts validated models for 3 stimuli, and provides a
31 new, experimentally testable model for a fourth (we will add suggested experiments to the paper). The main aim of
32 our paper is to publish our new hypotheses in order to stimulate multiple retina labs worldwide to tackle the difficult
33 neurophysiology experiments. In this manner, our theory could generate new experimental progress in future work.

34 **Simpler approaches do not suffice (R3)**: A single linear receptive field (RF) plus a nonlinearity (LN model) cannot
35 account for any of the 4 stimuli (indeed that is precisely why these stimuli are interesting). References from **R3** show
36 that ON/OFF pathways differ in their threshold as well as timing, and optimized two-pathway LN models could partially
37 capture the OSR [17] but *cannot* produce sufficient frequency-dependent shifting of the latency [18]. Thus the reported
38 asymmetries cannot produce the observed OSR response, and our new finding is that three pathway LN models can.

39 **Clarifying our contribution beyond previous work (R3)**: While building on a deep retina network from the authors
40 of [2], that work did not provide conceptual understanding of *how* the network generated responses to 4 highly structured
41 stimuli, and *whether* it generated those responses the same way the retina did. We provided such an explanation,
42 showing only 3 of 8 channels were required to generate responses to all 4 stimuli in an interpretable manner, thereby
43 demonstrating a single approach (natural scenes -> deep CNN -> model reduction) that can *simultaneously* discover
44 what was previously only discovered piecemeal across ≥ 10 papers. We feel this yields a major advance in providing a
45 "roadmap for neuroscience" (**R1**). Moreover, our method is primarily a *novel application* of attribution methods in
46 [9,10] to model reduction in neuroscience with validation in a biological circuit. From the NeurIPS call for papers, such
47 *application* papers are squarely within conference scope, and major advances in attribution methodology should not be
48 required for acceptance since that direction is orthogonal to our application to model reduction in neuroscience. We
49 will however revise to tone-down, discuss limitations, and clarify specific contributions (**R3,4**).

50 **Revising the text (R2)** We will follow **R2**'s excellent suggestions; we will shorten the intro, expand results, move info
51 from Fig. 2 caption to text, and provide more background on integrated gradients in the main, using the extra page for
52 the camera-ready. We note **R2** gave the lowest score (4 compared to 8 (**R1**) and 7 (**R4**)), despite being very positive
53 (**R2**: "highly original contribution with huge potential," "with the supplement, the methods are clear"). We hope, given
54 our restructuring, **R2** will be convinced that the revised version will be acceptable.

55 **Other comments (R1-R4)** Though we cannot address all remaining less major comments in the author response due
56 to lack of space, we assure reviewers we can easily do so in the revision. We are grateful for your excellent suggestions.