1 We thank the reviewers for their critical feedback. Quite honestly there were large gaps in the explanation and analysis
2 (as loosely clustered below: 10 of them), and we've made serious improvements to address nearly all:

3 **1. Missing actionable use cases (R2, R3); implications of "50.7% helped" unclear (R3)**. Our paper aims to increase
4 the understanding of neural networks in the same vein as insightful but not directly actionable papers like "Intriguing
5 properties of Neural Networks" (Christian Szegedy et al, 2013) and "Opening the Black Box..." (Shwartz-Ziv and
6 Tishby, 2017). We did find one use case (identify layers for freezing), but the primary goal of our paper is to provide a
7 new measurement device (LC) and to clarify and update our mental models of NN training by using it.

8 **2. Why is LC a good metric? Compare with FIM from [1,13] (R1)**. We believe both are useful. LC answers the
9 specific question of "how do we allocate change in loss?" and is beneficial vs FIM because it is *grounded to the loss* and
10 *signed*. *Grounded*: While [1] and [13] are illuminating, the FIM metric used may be rescaled arbitrarily (e.g. multiply
11 one relu layer weights by 2 and next by 0.5: loss stays the same but FIM of each layer changes and total FIM changes).
12 In contrast, LC is grounded, so its units are the same as loss (e.g. bits or nats for cross entropy) and its scale is fixed.
13 *Signed*: FIM is unsigned and thus could not yield our "50.7% help" nor "some layers move backward" conclusion. We
14 have clarified these points in the paper. *[See Update A in the figure below]*

15 **3. Justify use of batch gradient (R1)**. This part was poorly presented; we have updated the text *[Update B]*. To
16 measure training instead of training confounded with issues of memorization vs. generalization, we use the *train*
17 gradients instead of *val*. Observations like the last layer hurting are also more surprising on train vs. val. We use
18 *full-batch* gradients for analysis instead of single mini-batch gradient to measure learning in as noise-free a way as
19 possible. Note that the optimizer only ever sees mini-batch gradients, as is the usual case for SGD or Adam.

20 **4. Missing higher order terms, ignores curvature (R2); characterize approximation errors from first order (R1)**.
21 This part was confusingly presented; text updated to clarify *[Update C]*. Curvature is not ignored, as LC is approximated
22 with RK4, which is fourth-order accurate. First order was only mentioned as an example to illustrate the concept.
23 Approximation errors on all networks with both first order and RK4 have been added *[Update D]* (good idea!). In short,
24 first order can be off by a lot (it always overestimates the decrease in loss), but RK4 is within 0.3% total loss change.

25 **5. "Help" and "hurt" are too simplistic and myopic, loss may locally increase then decrease later (R2)**. Each LC
26 measure concerns the impact on loss at each specific iteration, and we agree that periods of hurting may be followed by
27 helping and vice versa. We see this often (e.g. see Fig 3d). By summing LC over time, we can take arbitrarily less
28 mypoic views (e.g. see Figs 4 and 5). We believe both microscopic and aggregate views can be useful.

29 **6. Last layer: different behavior and benefit of freezing already known (R2, R3)**. Thanks for these great references;
30 we have added a discussion *[Update E]*. We build on this previous work by proposing a theoretical explanation (phase-
31 delay) and validating the theory with careful experiments. See also: new multi-layer results in #7.

32 **7. How universally does last layer hurt? (R1)**. We don't claim this to be a universal phenomenon, but when it occurs
33 LC allows us to uncover it and take beneficial actions such as freezing layers. It is also not just a one-off observation: in
34 addition to ResNet, it also happens for AllCNN (mentioned in paper) where the last layer is conv, and we have now
35 measured it in VGG-based models as well, where the last 2-3 FC layers show the phenomenon *[Update G]*.

36 **8. Citations, oscillations not novel (R3)**. Thanks for these great references; we had not meant to claim discovery of
37 oscillation and have updated the text to clarify *[Update F]*. LC merely tracks such behavior on a per-parameter level
38 and can expose, for example, whether a parameter or layer is slowly making progress or the reverse over time.

39 **9. Intuition of results of freezing the first layer (R3)**. Indeed, given the observed behavior of freezing the first layer
40 (that other layers then help less), we should be careful not to over-interpret a positive per-layer LC as directly indicating
41 that that layer should be removed or frozen. The interaction effects between layers would seem to matter a lot, though
42 we do not fully understand these effects yet (though our phase-delay results are a step in the right direction).

43 **10. Learning is heavy tailed section not precise (R3)**. Thanks for the suggestion. We've updated it to a more precise
44 characterization: the LC distribution is closely modeled by a heavy-tailed Weibull distribution *[Update H, dashed line*
45 *shows the Weibull distribution, parameter $k = 0.53$ in example]*. This refines a detail in our mental model of training.



A, B, C. Clearer explanation of LC method and comparison with other methods

**Update Highlights**

D. Added table of approximation errors for RK4 (used in main findings) and first order

E. Citations for last layer behavior

F. Citations for oscillations

G. The last FC layers of two new VGG-styled networks also hurt

H. LC distribution closely fits a Weibull distribution (dashed line)