1 We thank the reviewers for their very useful comments.We are encouraged by reviewers thinking our work "is a novel
2 contribution" (R1), its topic "is very relevant to GLMs" (R2) and "seems like an important advance" (R3). We are glad
3 they thought our work has general significance in the neuroscience community at NeurIPS (R1) and many will find it
4 useful (R2). We address the reviewers' comments below and we will incorporate all the feedback. Rebuttal Figure 1
5 shows some selected panels (for space reasons) from Figs. 3 and 4 that we hope illustrate how changes proposed by the
6 reviewers will be incorporated. We will incorporate error bars by repeating the optimization procedure multiple times
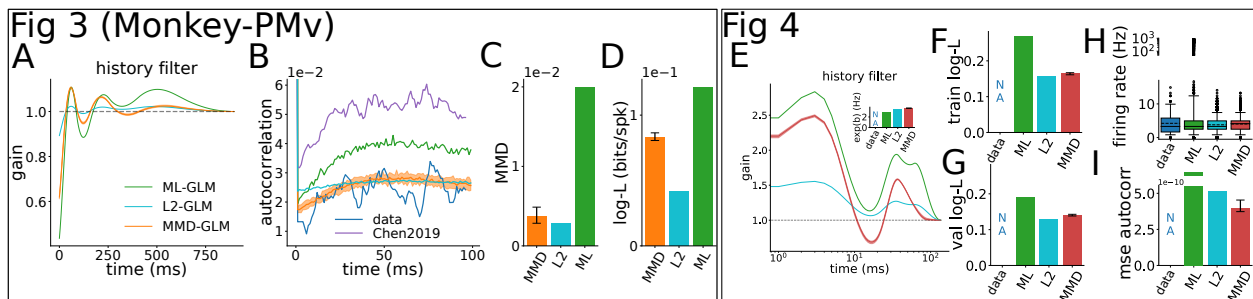7 and we will add L2 regularization on the history filter to compare with our method.

8 **R2 suggested regularizing the history filter.** As R2 noted, shrinkage of the history filter can lead to improved stability
9 of the model. Rebuttal Figure 1 shows results for L2 regularization on the history filter. Results suggest that increasing $\alpha$
10 in L2 regularization to match the data firing rate or stabilize the model (see next response for the choice of regularization
11 strength) in general leads to worse general performance than MMD regularization. In contrast, MMD regularization
12 can be used to explicitly match any quantity and not just reduce the size of the history filter coefficients as in L2
13 regularization. MMD regularization yielded parameters with higher likelihood values (Rebuttal Figure 1 D, F) that also
14 captured better the data autocorrelation (Rebuttal Figure 1 B, I) (data autocorrelation was smoothed to improve visuals).
15 Note that the (model-based) MMD in Rebuttal Figure 1C is smaller for L2 regularization than for MMD regularization,
16 but this is simply due to the L2 regularized model matching its free-running and data-conditioned distributions.

17 **Choosing regularizer weight ($\alpha$ in NLL + $\alpha$MMD) (R2).** There are multiple possible strategies for choosing the
18 hyperparameter $\alpha$—one could focus on matching one of the key statistics of the free-running model (e.g. mean firing
19 rate) with the data, or minimize the runaway excitation. For Figure 3 Monkey-PMv dataset we performed a grid search
20 and chose the smallest $\alpha$ value such that the mean firing rate of 8000 samples differed less than 10% from the mean
21 data firing rate. We used the same criteria to determine the regularizing weight in the added L2 regularization (Rebuttal
22 Figure 1 A-D). For Figure 3 Human Ctx dataset we chose the smallest $\alpha$ such that out of 2400 samples none showed a
23 diverging firing rate. For Rebuttal Figure 1 E-I we chose the smallest $\alpha$ such that out of 8000 samples none showed a
24 diverging firing rate. We would state in the final version how we choose $\alpha$ for all the models in Figure 4.

25 **Statistical significance and error bars (R1, R2).** As pointed out by R1and R2, our estimation procedure is stochastic.
26 Rebuttal Figure 1 shows panels from Figures 3 (Rebuttal Figure 1 A-D) and 4 (Rebuttal Figure 1 E-I) with the results of
27 repeating the optimization procedure 20 times. For Figure 4 panels, we illustrate variability for one choice of kernel
28 (same model-based kernel as in Figure 3, autocorrelation of the history filter convolved with spikes). The shading in
29 the curves and the error bars in the barplots represent the 95% percentiles of the distribution obtained by repeating
30 the optimization. The lines and bar heights represent the mean over optimization repetitions. We can see that for the
31 examples shown here the stochastic optimization with the proposed MLE initialization is robust and similar parameters
32 are found consistently. Similar robustness was observed when repeating the optimization procedure multiple times for
33 the Human Cortex dataset of Figure 3 (not shown but would be incorporated in the final version).

34 **Scaling (R3).** Our procedure is definitely more expensive than MLE as it involves sampling and computing MMD
35 during the optimization—in fact, we initialize the optimizer at the MLE. The time and space complexity varies with
36 many factors; e.g. score function or model-based MMD, biased or unbiased estimator and the choice of the kernel—for
37 some kernels, it is quadratic in the number of spikes, and for some other kernels, it only depends on the duration of
38 recording. All the optimizations in the submission take as much as a minute to five in a desktop computer (compared
39 with seconds for MLE) and there is room for improvement in the current implementation of the procedure. The authors
40 have experiences with designing computationally efficient kernel and optimization tricks, and we plan to continue to
41 improve the practicality of this novel concept.

42 **Generating from ML fit discarding run-away excitation trials (R2).** Removing "outlier" trials from the generative
43 model has several disadvantages: (1) it can be very inefficient, (2) ML cannot guarantee that the truncation will result in
44 a "good" distribution, and (3) interpretation of the model parameters is no longer straightforward.



Rebuttal Figure 1: **(A) Corresponds to 3B.. (B) Corresponds to 3C. (C, D) Correspond to 3D. (E) Corresponds to 4A. (F) Corresponds to 4B. (G) Corresponds to 4C. (H) Corresponds to 4D. (I) Corresponds to 4H.**