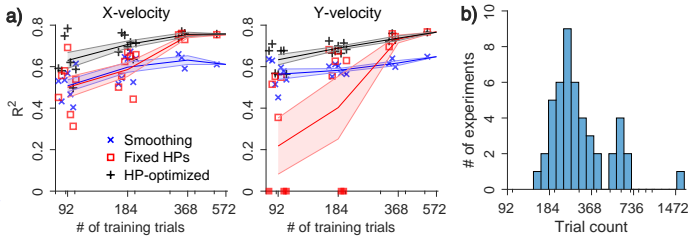1   **We appreciate the thoughtful feedback.** All reviewers noted that our *sample validation* (SV) and *coordinated dropout*
2   (CD) methods were novel with broad applicability. New analyses, clarifications, and proposed modifications are below.

3   **R1:** *Paper would be much stronger if ideas*
4   *were demonstrated on multiple real datasets* Done
5   (**Fig. 1a**). We used an open dataset [1] with a Ran-
6   dom Target task (different lab and experiment). We
7   found similar results to orig. Fig. 5, including the
8   range where HP opt helps, and the gap between op-
9   timized and fixed HPs. **R1:** *Description of typical*
10  *dataset sizes would help motivate the criticality of*
11  *the issue*; *Single small dataset is insufficient to estab-*
12  *lish general efficacy.* Agreed, we'll discuss. Typical



**Fig. 1.** **a**) Rand Targ task  **b**) # of trials for 47 experiments [1]

13  sizes largely vary, so for context we'll show the trial counts for 47 experiments from the open dataset (**Fig. 1b**; [1]).
14  These dataset sizes are typical, and many are in the range where HP opt is important. Note: our original dataset
15  (1836 trials) is actually *exceptionally large*, chosen so we could characterize HP opt vs. dataset size. **R1:** *Not clear*
16  *why "Monkey J Maze" is not used from the beginning... Synthetic data is unconvincing.* This is a key point. It is
17  important to clarify the necessity of tests on synthetic data, and may also help for readers without neural data experience.
18  **The synthetic data is critical - without it, it is very challenging to determine whether an approach results in**
19  **pathological overfitting**. Real neural data has no ground truth for direct comparison - there is no "true", measurable
20  firing rate. Common validation measures are problematic for detecting overfitting: *1)* Held-out likelihood of observed
21  data is somewhat noisy and requires assumptions. *2)* Decoding behavior, as we do, is a rough measure: only a small
22  fraction of neural activity correlates with behavior, and behavioral dynamics are quite slow. A precise characterization of
23  overfitting (orig. Fig. 1) and of the effectiveness of SV/CD (orig. Fig. 4) would be very challenging with real data. Since
24  SV & CD are the key innovations, we must thoroughly characterize them using data with a ground truth, and synthetic
25  data are the best option. To speed manuscript, we will move all synthetic data generation details to a supplement. **R1:**
26  *Existing regularization like denoising autoencoders (dAEs) should also be used as baselines. Motivation for completely*
27  *new techniques should be explained.* Great suggestion. We tested dAEs (**Fig. 2**),



**Fig. 2.** Denoising AE results

28  and motivation is now easily explained in the context of these results. We re-
29  peated orig. Fig. 1b using two common dAE approaches for discrete data: 'Zero
30  masking' and 'Salt and pepper noise' [2]. Important points: *1)* dAEs have a free
31  parameter (noise level). *2)* Depending on its setting, dAEs can still show patho-
32  logical overfitting. *3)* Some settings can even reduce performance. *4)* It is not
33  possible to know how to set dAE noise *a priori*. Our methods bypass these limi-
34  tations (see orig. Fig. 4), providing a reliable metric to measure (SV) or completely block (CD) pathological overfitting.
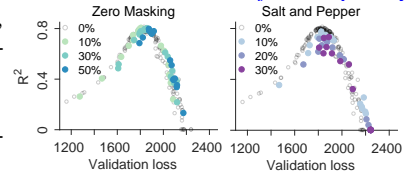35  **R2:** *Discuss if method can be extended to other data sets.* Good point, will add. Techniques should be applicable when
36  forecasting time series from sparse data, especially when HP or architecture searches are important. Examples are usage
37  at electrical vehicle charging stations, taxi/rideshare calls, etc.. We're currently trying to apply this to generative models
38  for LIDAR/RADAR data for autonomous cars (e.g., following [3]). **R3:** *Would raise my score with the inclusion of some*
39  *details that were missing... complete formulation of the generative model and inference procedure.* Good suggestion. We
40  will add this information. R3's description of objective was accurate. **R3:** *State validation loss and how it is computed...*
41  *Useful to fully describe LFADS model, at least in appendix.* Apologies for omissions, will add. **R3:** *Does the model still*
42  *exhibit pathological overfitting with AR prior included?* Yes, and we were surprised by this (all the results in paper
43  are with AR prior included). Key problem is AR prior is learnable, and model can adapt it to get better predictions by
44  overfitting to spikes via inputs. Forcing a minimum AR prior autocorrelation might prevent overfitting, but might also
45  prevent the model from capturing rapid changes. **R3:** *What HP settings provided "good" fits? Would be interesting to*
46  *include a discussion, including how this might vary across dataset size.* Agreed, including settings/ranges will be helpful.
47  Further, these methods enabled dynamic HP opt (changing HPs during training) using population based training [4].
48  This somewhat surprisingly yields even higher performance by learning schedules for different HPs (e.g., KL penalty
49  is set high during early training, but decreases over time). We'll add this discussion. **R3:** *Is full-split CD necessary,*
50  *or could you also split the data into input only, shared, and output only splits?* This is very interesting, we've been

51  thinking about this also. The proposed 'Partial CD' approach might help when observed number
52  of neurons is similar to the underlying dimensionality, and fully splitting data via CD may limit
53  training. Without Full CD, though, a method is needed to detect/prevent overfitting. SV fills this
54  role. As suggested, we turn CD on, and then allow some fraction of the data (searchable HP) to
55  be shared as input and output. Preliminary tests on small sets of randomly drawn neurons (Monkey
56  J Maze data, 25 per draw) show promising results: Partial CD outperforms Full CD in 8/10 models



**Fig. 3.** Partial CD

57  tested. Thorough tests will help delineate conditions where Partial CD helps.

58  [1] J E O'Doherty et al. http://doi.org/10.5281/zenodo.583331, 2017.    [2] P Vincent et al. *J. Mach. Learn. Res.*,
59  11:3371–3408, 2010.    [3] L Caccia et al. *arXiv:1812.01180*.    [4] M Jaderberg et al. *arXiv:1711.09846*, 2017.