

1 We thank all reviewers for their thoughtful feedback and suggestions. We are encouraged that all reviewers found the
2 paper well-written and agreed that the technical ideas are important and novel. Based on suggestions from the reviewers,
3 we have further improved the clarity of the paper and included additional details.

4 **(R1) - Interaction between HiPPO and cost function.** The overall HiPPO-RNN model (Sec. 4) uses the same loss
5 function as the baselines (e.g. cross-entropy for classification, as in pMNIST). Note that the HiPPO-RNN cell in Fig.
6 2 depicts only the recurrent dynamics; as is convention, each cell also outputs the state h_t for another layer (e.g. softmax)
7 to compute the final prediction. When incorporated in an end-to-end model such as an RNN, the HiPPO module (blue
8 circle in Fig. 2) can simply be viewed as a fixed, non-trainable component – although it has the additional interpretation
9 of memorizing the history of features (in L2 space). Since it only affects feature representation, the HiPPO component
10 does not impact the choice of output space and loss function, which is inherited from the base RNN. These details will be
11 clarified in the next version of the paper, and are also explicit in the submitted supplementary code which will be released.

12 **(R1) - Reproducibility and experiment details.** Additional space permitted, we have moved much of the architecture
13 details in Appendix F.1 to the main body, including (*) a better diagram of the HiPPO-RNN cell, (*) a discussion of
14 the output layer and loss function, (*) a comparison of network architecture, size, and hyperparameters against all RNN
15 baselines (in summary, they are controlled to be all roughly equal), (*) loss function details for all tasks in Appendix F.

16 **(R1) - Technical details.** We appreciate R1’s close reading of the paper and suggestions for improving the technical
17 presentation. We agree with all details on integrable functions, the Leibniz rule, the Copying task hyperparameter, and
18 the interpretation of coefficients $c(t)$. Thanks for the suggestions and corrections!

19 **(R1, R2) - Finding the polynomial order N .** In the setting of approximating smooth functions, N can be calculated from
20 the explicit error bounds (Sec. 3, Prop. 6). In our end-to-end experiments, N is a hyperparameter analogous to the hidden
21 dimension of standard RNN/LSTM models; we set N equal to the hidden dim. d of the baselines, so that all HiPPO methods
22 have the same number of hidden units and the same (or smaller) model size as the LSTM baseline (line 1094-1104).

23 **(R2) - HiPPO, RNNs, and universal approximators for dynamical systems.** Similar to universal approximation
24 theorems, HiPPO formalizes the ability of RNNs to approximate functions. HiPPO shows that optimal compression of
25 signals is governed by linear recurrent systems (Thm. 1, 2). Unlike universal approximation theorems, HiPPO shows
26 non-asymptotic error rates on the approximation in terms of memory capacity (Prop. 6). The HiPPO framework is a bridge
27 between RNNs and dynamical systems, and explains how to derive the form of a modern RNN update from first principles.

28 **(R2) - Directly fitting timeseries in the input space.** HiPPO can be used directly on input features, which is how the
29 basic HiPPO framework is described. For example, the function approximation experiments (Sec. 4.3 and App. F.8)
30 can be interpreted as fitting a long input timeseries with a budget of 256 features. We remark that the presented figures
31 (Fig. 3, 9, 10) are a good visual indicator of the approximation dynamics of HiPPO.

32 **(R2, R4) - Short-term or more local time series tasks.** While HiPPO was designed to capture long-range dependencies,
33 several of our experiments involve tasks requiring shorter memory: trajectory classification in Sec. 4.2 (length ≈ 100),
34 chaotic dynamics simulation in App. F.7 (length 15), and the IMDB dataset in App. F.6 (some inputs of length < 100).

35 **(R3) - Experiments.** R3’s only concern is that “the paper needs more experiments to better support the proposed model...
36 such as speech processing, NLP, etc.” Although we acknowledge R3’s suggestion, we note that the submission already
37 builds up from the foundational theory to a comprehensive range of experiments including an NLP application. Indeed, in
38 addition to directly validating the proposed model and online function approximation theory on synthetic experiments, we
39 have included experiments on memory benchmarks, image classification, trajectory classification, timescale robustness,
40 chaotic dynamics prediction, and NLP (the IMDB movie review benchmark); some of these are in Appendix F due to
41 space constraints. We are excited about pursuing further applications in follow-up work.

42 **(R3) - Baselines.** Although R3 asks about more memory RNN baselines, the memory experiments in the paper include a
43 variety of well-known memory models (including RNNs), chosen for their relevance to each task. Our permuted MNIST re-
44 sults explicitly compare against 15 sequence model baselines (Tab. 4, 5) to establish SoTA; this also establishes a large num-
45 ber of implicit comparisons, e.g., the LMU outperformed 8 other memory RNN models in (Voelker 2019). We ran an addi-
46 tional baseline suggested by R3, the MinimalRNN, on the memory benchmarks (Sec. 4.1). This model failed to solve the
47 Copying task and achieved 89.1 accuracy on the permuted MNIST benchmark, compared to the SoTA 98.3% of our model.

48 **(R4) - Why LegS is faster than LMU.** The speed difference in Sec. 4.3 is explained by our faster algorithm for LegS.
49 Because of the approximation made by LMU due to its measure (Sec. 2.3 lines 158-161), the LMU transition matrix
50 A (Thm. 1) is not triangular, in contrast to that of LegS (Thm. 2). Numerically stable discretizations (App. B.3) require
51 inverting this matrix, which is efficient in the LegS case (Prop. 4, proved in App. E.2) but not known for the LMU case.

52 **(R4) - Performance depends on measure and task.** We fully agree with R4’s observation that some tasks are more
53 suited to LegS due to the uniform memorization prior. In fact, rather than seeing this as a weakness, we argue that
54 a primary contribution of this paper is introducing the technical framework that exposes these very tradeoffs. This
55 framework allowed us to (1) explain the memory mechanism of existing methods (LMU, LSTM, Fourier Recurrent
56 Unit, etc.) in terms of their approximation measures, (2) introduce new methods that may be more appropriate in different
57 settings such as very long memory and mis-specified timescales, and (3) theoretically analyze and contrast different
58 approaches based on their underlying mechanisms. We thank R4 for bringing up this insightful point!