1 Thanks to the reviewers for their valuable feedback. We are encouraged by their endorsements that our MomentumRNN
2 framework: 1) makes a very novel and thought-provoking [**R3**] connection between RNN and optimization [**R3**,**R4**],
3 2) provides a thorough analysis integrating various forms of momentum into different RNN types [**R2**], and 3) leads to
4 decent improvements over baselines [**R1**,**R2**,**R3**]. Below we address the concerns raised by the reviewers.

5 [**R1**,**R4**] **Momentum cell is similar to gated update** $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{h}_t$ **of LSTMs and residual mappings.** We
6 respectfully disagree with this point. We believe there is a misunderstanding. The MomentumRNN cell updates as:
7 $\mathbf{v}_t = \mu\mathbf{v}_{t-1} + s\mathbf{W}\mathbf{x}_t; \mathbf{h}_t = \sigma(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{v}_t)$, which introduces an auxiliary state $\mathbf{v}_t$, inspired by Nesterov momentum.
8 First, for the hidden state dynamics, the two-step update in momentum cell is equivalent to $\mathbf{h}_t = \sigma(\mathbf{U}(\mathbf{h}_{t-1} -$
9 $\mu\mathbf{h}_{t-2}) + \mu\sigma^{-1}(\mathbf{h}_{t-1}) + s\mathbf{W}\mathbf{x}_t)$, i.e., $\mathbf{h}_t$ directly depends on $\mathbf{h}_{t-1}$ and $\mathbf{h}_{t-2}$; in LSTM, $\mathbf{h}_t$ only directly depends on $\mathbf{h}_{t-1}$.
10 $\mu\sigma^{-1}(\mathbf{h}_{t-1})$ is the key term that helps alleviate the vanishing gradient. Second, $\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{h}_t$ differs from
11 $\mathbf{v}_t = \mu\mathbf{v}_{t-1} + s\mathbf{W}\mathbf{x}_t$, where $\mathbf{f}$ and $\mathbf{i}$ are tensors while $\mu$ and $s$ are scalars; also, $\mathbf{h}_t$ is the hidden state while $\mathbf{x}_t$ is the
12 input. Third, $\mathbf{v}_t = \mu\mathbf{v}_{t-1} + s\mathbf{W}\mathbf{x}_t$ differs from residual mapping ($\mathbf{v}_t = \mathbf{v}_{t-1} + \mathcal{F}(\mathbf{v}_{t-1})$). As pointed out by [**R3**], "*In*
13 *LSTM, the cell state performs additive integration of the input so that the gradients do not vanish. However, LSTM's*
14 *sigmoid and tanh gates make gradients vanish. MomentumRNN performs additive integration of the input without any*
15 *additional gates and hence has a much better gradient flow.*" We hope reviewers can reevaluate this crucial point.

16 [**R1**] **Analogy between MomentumRNN and optimization methods.** We do not aim to equate RNNs with gradient
17 descent. Rather, we aim to bring in new ideas from optimization to design better RNNs. The nonlinear activation can be
18 considered as a projection, and, based on our experiments, the acceleration from momentum is preserved. Moreover, the
19 resulting momentum cell derived from optimization algorithms with momentum can alleviate the vanishing gradient.

20 [**R1**] **Analysis is w.r.t. the vanilla RNN. Also, it is not obvious that an appropriate** $\mu$ **exists.** In Sec. 2.3, we proved
21 that $\mu\sigma^{-1}(\mathbf{h}_{t-1})$ in the MomentumRNN cell $\mathbf{h}_t = \sigma(\mathbf{U}(\mathbf{h}_{t-1} - \mu\mathbf{h}_{t-2}) + \mu\sigma^{-1}(\mathbf{h}_{t-1}) + s\mathbf{W}\mathbf{x}_t)$ is key to alleviate the
22 vanishing gradient in vanilla RNNs, which generalizes to the analysis of other RNNs. Our empirical study in Fig. 2 is not
23 for a vanilla RNN, but for the SOTA DTRIV model (see Appendix A.4). Our experiments show that MomentumRNNs
24 significantly outperform all studied baseline RNN models, which empirically validates that an appropriate $\mu$ exists.

25 [**R1**,**R4**] **Toy datasets and lack of comparison to other SOTA RNN baselines.** We used the (P)MNIST, TIMIT, and
26 Penn TreeBank (PTB) benchmarks, which are not toy. Our momentum-based method can be applied to many other
27 RNNs, including those that that address vanishing gradients, to improve their performance. In Sec. 3 and 4, we have
28 shown this in the case of LSTM and expRNN (DTRIV), a SOTA RNN model. Our approach outperforms both.

29 [**R1**,**R3**] **More related work on long-term dependencies and saturations of LSTM.** We appreciate the suggestions
30 and will add JANET, NRU, and more papers on alleviating long terms dependency issues in the revision.

31 [**R2**] **Why the MomentumRNN provides improvements over RNN is not clear.** MomentumRNN mitigates vanish-
32 ing gradients as analyzed in Sec 2.3. The reason it can theoretically accelerate convergence, improve robustness, and
33 lead to higher performance is under our study. In particular, we are studying the continuum limit of the MomentumRNN.

34 [**R2**] **Compare the computation cost of MomentumRNN & RNN to reach similar acc.** When training on the
35 PMNIST task using 256 hidden units, we observe that *to reach 92.29% test acc. for LSTM (see Tab. 1), LSTM needs*
36 **699min** *while MomentumLSTM & RMSPropLSTM (our best model for this task) only need* **416min** *and* **403min***, resp.*

37 [**R3**] **Why should RMSPropLSTM work?** RMSPropLSTM inherits its adaptive step size from RMSProp. In training,
38 $\mathbf{W}\mathbf{x}_t$ is rescaled adaptively and can improve training and enhance performance for certain tasks.

39 [**R3**] **Why using different optimizers for different tasks? Results on using Adam for all models.** We used the
40 default optimizers that achieve SOTA results for different tasks, e.g., we used *RMSProp for (P)MNIST, Adam for TIMIT,*
41 *and SGD for PTB*. Using Adam to train LSTM, MomentumLSTM, AdamLSTM, RMSPropLSTM, and SRLSTM with
42 256 hidden units on PMNIST, we obtain test accuracy of $92.05 \pm 0.63\%$, $92.47 \pm 0.35\%$, $92.53 \pm 0.26\%$, $93.86 \pm 0.24\%$
43 and $92.17 \pm 1.37\%$, resp. Adam trainings for DTRIV and MomentumDTRIV with 512 hidden units on PMNIST yield
44 the test accuracy of $95.78 \pm 0.21\%$ and $96.01 \pm 0.10\%$, resp. Adam training yields worse results than SGD on PTB.

45 [**R3**] **1) Why using only the MomentumRNN in language modelling task? 2) Fig. 6 shows that often momen-**
46 **tum=0 is the best thing to do? 3) Why is forget gate initialized to -4? 1)** We compare 3-layer MomentumLSTM
47 with 3-layer LSTM in our paper (see Tab. 3). The test PPL for AdamLSTM, RMSPropLSTM, and SRLSTM are
48 $61.11 \pm 0.31$, $64.53 \pm 0.20$, and $58.83 \pm 0.62$, resp. **2)** In Fig. 6, momentum=0 yields the best results only for language
49 modeling task. We believe that if we do finer-scale search for momentum and step size, we can obtain a better result with
50 a non-zero momentum. **3)** For the TIMIT task, initializing the forget gate bias to -4 is suggested in [arXiv:1707.09520].

51 [**R2**,**R4**] **Additional hyper-parameters & their sensitivity.** We search $\mu$ & $s$ independently for a few options which
52 is not very sensitive and does not increase the computational cost much. Also, trainable parameters is under our study.

53 [**R4**] **Copy/adding tasks.** In copying (adding) task for sequences of length 2K (750), our MomentumLSTM, AdamL-
54 STM, RMSPropLSTM, and SRLSTM achieve the final training loss of 0.009 (0.162), **0.004** (0.006), 0.008 (**0.004**), and 0.01
55 (0.166), resp. while the training loss of the baseline LSTM is 0.01 (0.162). AdamLSTM and RMSPropLSTM converge
56 faster than LSTM in both tasks. Moreover, AdamDTRIV and RMSPropDTRIV converge to the same final training loss
57 remarkably faster than DTRIV in both tasks. Detailed results are included in the revision.