

1 **Response to Reviewer 1:**

2 **"It is not clear what terms eventually vanish and why."** Terms multiplied by weight matrix W (red boxes) in
3 equation after line 97 vanish if the largest eigenvalue of matrix W is less than 1. Multiplying it repeatedly causes the
4 vanishing gradient. This reduces equation after line 97 to equation in line 100.

5 **"Euclidean distance is penalized if saliency goes over 1."** We use the normalized Euclidean distance as mentioned
6 in line 157 so the distance will never go over 1.

7 **"The fMRI dataset do not show anything more than what previously laid out"** The goal of the dataset is to show
8 a real-case application where we are interested in seeing how features importance change across time. Previous work
9 [1] that has been done on this dataset gives a single interpretation for the entire time series and does not show how
10 features change across time. To clarify our point we will add the off-task time to our experiments (time where subject is
11 listening to instructions not performing the actual task) and we will investigate the ability of our model to ignore this
12 time and put the importance on the on-task time.

13 **"The definition of self-attention is not up-to-date."** We choose to use the definition of self-attention of Lin et al.
14 (2017) since they are applying self-attention to a RNN, similar to our case. Also, we do not use the same attention
15 function described by Vaswani et al. (2017).

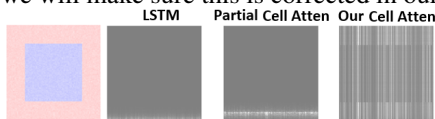
16 **"Improved saliency comes from the attention done in the last time steps of the RNN."** We disagree. If this was
17 true then removing attention from first time steps would make no difference which is not the case. Figure 1a shows an
18 experiment where we applied attention only to the last 10 time steps (referred to as partial cell attention) for middle
19 box dataset, saliency still vanishes. Our proposed method improves saliency because, at each time step, cell-attention
20 attends to different inputs from current or previous time steps preserving importance through time.

21 **"Proposed approach ignores the temporal nature of the problem"**. We disagree. Our entire paper is based on time
22 series, permuting data in time and producing different saliency is the entire purpose of our synthetic dataset. The
23 moving box experiment in line 158 and figure 1 in supplementary material shows a clear example where the only
24 difference between samples is their location in time, it is very clear from the experiments that we do *NOT* ignore the
25 temporal nature of data.

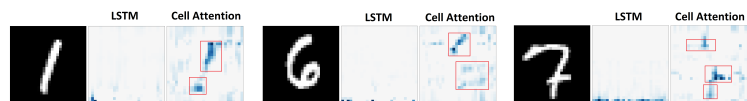
26 **Minor comments.** For weighted Jaccard, we compare the saliency with the absolute value of synthetic data sample, we
27 will update this in the final version. We will add numbering to all equations, correct mentioned typos and fix coloring
28 for figure 1 in supplementary material.

29 **Response to Reviewer 2:** Thank you for your comments. In the original version of the paper, we mentioned related
30 work briefly in the introduction but we did not have an entire section dedicated to related work due to space limitation.
31 In the revised draft, we will add a related work section and make sure we cite and explain all papers you listed in this
32 section along with others. Our scope in this paper is on studying saliency of RNNs where we propose an approach to
33 resolve the vanishing saliency that hinders the interpretation of such networks. Thus, we have kept the comparison with
34 vanilla attention and other non-recurrent network architectures to our future work. Thank you for suggesting to add
35 other standard benchmarks to our experiments. Upon your suggestion, we decided to add two new benchmarks. (1) We
36 use MNIST dataset as a time series data where one dimension of the 2D images acts as the time axis (a 28×28 image
37 is turned into a sequence of 28 time steps, each of which is a vector of 28 features). We choose MNIST because it offers
38 an interpretable visualization. Figure 1b is an example of a saliency map produced for vanilla LSTM and our proposed
39 LSTM + cell attention. (2) CMU Multimodal Opinion Sentiment Intensity (MOSI, Zadeh et. al 2016), a dataset of
40 opinion level sentiment intensity in online videos. In the final version, we will include these new experimental results.
41 Finally, we will include distributions of weight matrices of the network, as suggested.

42 **Response to Reviewer 3:** Thank you for your comments. You are correct that the saliency is computed for each input
43 captured and accumulated till the current time step. We will make sure to make this point more clear in the final
44 manuscript. The accumulation effect is reduced by the approximation mentioned in the paragraph under line 121. We
45 called our method cell attention because its attention is on the cell level rather than hidden layer level although we
46 understand your concern about this name and how this might create some confusion. We may consider changing the
47 name of the method to Recurrent Attention. For lines 118 and 119: A_t has dimensions $r \times t$ where t is the number of
48 time steps in the current input, A_t has a weight for each time step; weight of all time steps sum up to 1. For lines 120
49 and 121 M is flattened to a vector of length $r * N$ and W_M is a matrix of $h \times (r * N)$. Thank you for pointing this out
50 we will make sure this is corrected in our final version.



(a) Experiment showing that having cell-attention **only** in the last time steps (Partial Cell Atten.) still produces vanishing saliency.



(b) An example of saliency map produced for MNIST, when treated as a time series. Saliency vanishes for vanilla LSTM while our proposed model is able to detect important features.