

1 We thank the reviewers for insightful remarks and comments that help to considerably improve our manuscript. We  
2 address the most important ones in detail below. Before doing so, we highlight a comment from R3 in order to make an  
3 important clarification about the scope of our contribution. *“It is well known that an attention mechanism would reduce  
4 gradient vanishing. It feels trivial to me as there is a direct connection for gradients to pass. [...] it’s not hard to see that  
5 having fewer things in the softmax (i.e. sparse attention) would increase the attention weights, thus improved gradient  
6 flow.”* We are in complete agreement and recognize that the very mechanism of (self-)attention was designed to improve  
7 gradient propagation over long sequences, and that sparsity is a good way to keep complexity costs low. However, to the  
8 best of our knowledge, there is currently limited understanding of gradient scaling properties in the presence of attention.  
9 Much like work from the ’90s established formal results for gradient exploding/vanishing in deep/recurrent networks, we  
10 believe it is crucial to establish similar theoretical tools for attention mechanisms, as these methods are under intense  
11 development where scalability and complexity are important issues. Our main aim is to contribute to this direction  
12 with a thorough analysis of gradient propagation in self-attentive systems, outlined in clear theorem statements which  
13 precisely quantify the intuition expressed by R3. These results are not trivial (see proofs in appendices), and offer valuable  
14 guarantees for attention mechanism development. The proposed relevancy mechanism and accompanying experiments,  
15 building on established work, are meant to illustrate how our theorems can be concretely exploited. We chose simple  
16 tasks for their ease of interpretation, and their variety of computational demands (memorization, prediction, RL, etc.).  
17 As is clearly indicated in the text, it is not our goal to propose this method “as is” in a race for state-of-the-art. Rather,  
18 we think it achieves its goal of providing a firm theoretical footing for a wide class of self-attention methods as well as  
19 suggesting future direction of methods based on concrete scaling guarantees and a balance between sparsity and attention  
20 complexity. We recognize that reviewers may have based their evaluation as they would have in a method paper, and we  
21 kindly invite them to reconsider the value of our experiments in the broader context of our theoretical contributions. We  
22 also thank reviewers for their additional minor comments not explicitly addressed here and agree to implement them.

23 **R1: Q** *“The authors didn’t spell out the relation between  $\kappa$  and  $d$ : higher  $\kappa$  tends to have smaller  $d$ . This relation is  
24 discussed in the paragraph of line-173 but not reflected in the formula of theorem-2.”* **A:** We thank R1 for this important  
25 remark. As it stands, Theorem 2 offers scaling relationship for any  $\kappa$  and  $d$ . However in practice,  $\kappa$  and  $d$  co-vary in ways  
26 that depend on the task’s underlying relevancy structure, a point that is explained in detail in Appendix C (see Fig 3)  
27 which explores trade-offs between  $\rho$  ( $= \kappa - \nu$ ) and gradient propagation (implicitly depending on  $d$ ). We will sharpen  
28 this discussion in the main text. **Q** *“In experiments, the authors mentioned the proposed model is “faster” to train but  
29 didn’t give any quantitative results.”* **A:** We will clarify this in the text and specifically point to Fig 4 in the Appendix D  
30 (or move it to the main text space permitting) which shows ReLSTM/ReLNN learns the copy and denoise tasks with  
31 significantly fewer number of updates as compared to other baselines.

32 **R2: Q** *“Line 145, how can Theorem 1 be related to the early attention mechanism [1]? As the attention weights are  
33 computed adaptively, it is unlikely that they are uniform.”* **A:** We recognize that uniform attention weights are unlikely  
34 in practice. This theorem offers a form of “worst case” guarantees which is applicable in practice for two reasons. (1)  
35 Typically, attention weights are initialized uniformly. (2) We experimentally verified that gradient propagation remains  
36 stable throughout training for a fully self-attentive RNN, see Fig 2 (Section 6). We will further clarify this in the text.  
37 **Q** *“What is the advantage of the proposed method over MANNs? / how are MANNs related to the Theorem 2? How are  
38 the paper’s theoretical findings different from [5], wherein gradient norm of self-attentive RNN is also quantified?”* **A:** We  
39 thank R2 for this insightful remark and acknowledge the lack of discussion about the MANN model class. Most instances  
40 operate on involved memory addressing/retrieval mechanisms that keep complexity costs low but, to our knowledge, do not  
41 offer gradient propagation guarantees like Theorem 2. Some instances, such as that proposed in [5], offer more similarities  
42 to our approach, and future adaptations of our Theorem 2 is an exciting future direction. Our contribution complements  
43 that of [5] in two subtle but important ways. (1) The proof of Theorem 2 in [5] describes an approximation of information  
44 propagation via recurrence, with “skip connections” contributions accounted only once in isolation. In self-attention,  
45 gradient propagates via a mix of skip and recurrent connections, leading to multiple gradient paths. Enumerating these  
46 paths and incorporating their contributions to gradient norms is the crux of our results, yielding complete and precise  
47 expansions of gradient terms applicable to a wide range of models, including Transformers. (2) The proposed method  
48 in [5] is an optimal memory writing schedule. In contrast, our method relies on relevancy for memory writing, allowing  
49 more flexibility in complexity scaling. We will incorporate this discussion in the revised text.

50 **R3: Q** *“augmenting RNNs with sparse attention to prevent gradient vanishing is not novel in itself [19]. The only novel  
51 part of the model is relevancy screening, but why that approach is better than other sparsity methods has no grounding in  
52 the theoretical analysis in the paper [...]”* **A:** We agree with R3 that sparse attention is not novel. However, quantifying  
53 how much sparsity contributes to gradient norm is new. The theorems we provide borrow notation from [19], but apply  
54 to a large class of self-attentive mechanisms as well. To see why theorem 2 directly applies to the relevancy screening  
55 mechanism, it suffices to take  $\kappa = \rho + \nu$ , as at each time step we are attending to the  $\rho$  states from the relevant set  
56 and the  $\nu$  states from the short term buffer. To highlight how our relevancy screening is grounded in the theory, see Fig  
57 3 in Appendix C, where we perform an experimental trade-off analysis between  $\kappa$  and  $d$  by tweaking  $\rho$  and  $\nu$  for our  
58 relevancy screening mechanism. This will be clarified in the text. We would like to emphasize that theorem 2 applies  
59 to any  $\kappa$ -sparse self-attention recurrent model, but we harness the structural understanding derived from the theoretical  
60 framework to propose a more scalable form of sparse self-attentive RNN.

61 **Q** *“Experiments: I think the experiments in the paper are quite weak [...]”* **A:** We appreciate this remark and refer to  
62 this rebuttal’s opening statement about scope. We also point out that our methods appreciably improves generalization  
63 to out-of-distribution samples over baseline. This is a promising avenue and will be further discussed in the text. We also  
64 want to thank R3 for pointing out the Transformer-XL reference, which we will make sure to include in the main text.