

1 **Review 1:** We thank the reviewer for the comments and suggestions. We will fix all typos pointed out in the review.
2 Below we address the two main issues raised by the reviewer .

3 **How is “sum” harder than “unique sum”?** Indeed accuracy on the *unique-sum* task is better than those of *sum*. The
4 reason is that for unique-sum the input domain is $\{0, \dots, 9\}$, and for sum it is $\{0, \dots, 99\}$ (we followed the setup in the
5 Janossy paper). This makes it possible to learn unique-sum by simply having a 10 dimensional state where dimension i
6 is “on” when digit i appears. For the rebuttal, we trained unique-sum with $\{0, \dots, 99\}$ and accuracy was much lower
7 than sum.

8 **How can DeepSets fail for $n = 15$?** DeepSets failure here is due to convergence to a bad local minimum. We used the
9 standard optimization setup for DeepSets (in the Janossy code). We believe that one advantage of SIRE is improved
10 robustness to optimization failures, because the regularizer also makes optimization easier by using shorter sequences.

11 **Baseline results for “half range”.** For the rebuttal, we ran the *half-range* experiment with DeepSets and it failed to
12 converge to a good solution, resulting in a model which outputs predictions based on the statistics of the data. The
13 resulting accuracy results numbers are (std in brackets): $0.050(4 \cdot 10^{-6})$, $0.077(1.4 \cdot 10^{-4})$, $0.0895(2.5 \cdot 10^{-7})$ for
14 $n = 10, 15, 20$. These are much worse than what SIRE obtains.

15 **Review 2:** We thank the reviewer for the comments. We will clarify definitions 2 and 3 according to the suggestions.
16 **Scalability of the proposed method.** The main focus of our work is to highlight the representational advantages of
17 RNNs, which makes them a very relevant part of the toolbox in invariant modeling. While this comes at a certain
18 computational cost, we believe there are many cases where the benefit in accuracy will make this attractive in practice
19 as demonstrated theoretically and empirically throughout the paper. A similar trade-off is evident in natural language
20 processing where recurrent methods such as GPT are widely used though there exists competing non-recurrent language
21 models such as BERT.

22 **Review 3:** We thank the reviewer for the through review and insights. We address each point raised below.

23 **How is SIRE sampled?** We estimate R_{SIRE} by applying a random permutation over input sequences x_1, \dots, x_n
24 followed by truncating the permuted input sequence to result in two sequences: $x'_1, \dots, x'_{k-1}, x'_k$ and $x'_1, \dots, x'_k, x'_{k-1}$
25 (i.e., the two last elements are switched). This simple scheme worked well, but we agree that other distributions can be
26 explored. Also, note that when the data is generated from S_n the above scheme is uniform over $\mathcal{S}_{\mathcal{D}, \Theta}$.

27 **What is the complexity of sampling s in estimating SIRE?** As mentioned above, SIRE is estimated by sampling
28 sequence pairs. The sampling operation is linear in sequence length, and then each one of the RNN is applied to the two
29 sequences. Thus, cost is similar to regular training.

30 **Number of SIRE Samples.** For each original input sequence, we sample just one pair of sequences for the SIRE
31 regularizer. Of course, for a given sequence, different pairs may be sampled each time this sequence is used.

32 **Def. 3, variables i and t .** Indeed, the definition requires $2 \leq t \in \mathbb{N}$ and $i < t$. We will clarify this.

33 **In Theorem 5, it should be “if and only if”?** The condition in the theorem is slightly stronger than subset-permutation-
34 invariance, because the sequence that results in adding x_1, x_2 may not correspond to a permutation of an actual input
35 sequence. The theorem can be easily adapted to make it “iff”, but notation would be more cumbersome.

36 **What kind of “semi”-invariance is compatible with SIRE?** In problems that are not permutation invariance, SIRE
37 will seek RNNs that are “as close as possible” to permutation invariant, where closeness is measured by SIRE. It will
38 indeed be interesting to further analyze this tradeoff, and we leave this for future work.

39 **How far can one push the sequence length until it no longer works?** The answer to this question depends heavily
40 on the task at hand and is mostly a matter of optimization. For the sum task, performance begins to deteriorate on
41 relatively short sequences - for $n = 15$ the average accuracy of an RNN using SIRE regularization is 0.957 (0.022). For
42 the purpose of the rebuttal we ran the same experiment with $n = 17$ and got an accuracy of 0.913 (0.003).

43 **Review 4:** The authors wish to thank the reviewer for the comments and points raised. We address the concerns below.

44 **Can the proposed method be used for permutation equivariance?** Yes. In a multi-output setting, we can add a
45 regularizer similar to SIRE, but that requires that adding x_1, x_2 to state s will generate outputs y_1, y_2 whereas adding
46 x_2, x_1 will generate outputs y_2, y_1 . This will result in similar guarantees to what we have, but for equivariance.

47 **Alternative aggregation methods.** This is indeed a main motivation for our work. In DeepSets and related
48 methods, the choice of aggregation method is key to obtaining good performance. While in principle one can check
49 $\{min, max, sum\}$ there are many other possibilities. On the other hand, in our case, the aggregation function is
50 learned automatically by the RNN, and as we show, the RNN works for cases corresponding to different aggregations.

51 **Sequential computation.** Recurrent architectures are used extensively in many machine learning applications, from
52 generative models in NLP (e.g., the GPT model is state of the art in language generation, and arguably outperforms
53 BERT which is non-recurrent) to video processing (where 3D convolutions and RNN over 2D convolutions are both
54 commonly used). The reviewer is correct that their sequential nature makes them less amenable to GPU speedups.
55 However, their representational power compensates for this limitation in many cases. Our argument in this paper is that
56 RNNs are very relevant for permutation invariant modeling, and we demonstrate this empirically in several cases.

57