

1 We thank the reviewers for their insightful comments. We first clarify our approach and then address specific concerns.

2 **R1, R2 *Forward-backward asymmetry and decoding strategy.*** NAOMI efficiently uses forward and backward hidden
3 states (h^f, h^b). Note that encoder and decoder share weights. For example, consider a situation where only x_0 and
4 x_8 are known, and we wish to impute $x_{1 \rightarrow 7}$ (x_1 to x_7). We first predict the mid-point $\hat{x}_4 = g(h_0^f, h_8^b)$ and update the
5 backward hidden states $h_{8 \rightarrow 4}^b$. Given \hat{x}_4 , we predict $\hat{x}_2 = g(h_0^f, h_4^b)$ and update $h_{4 \rightarrow 2}^b$. Recursively, given x_0, \hat{x}_2 , we
6 predict \hat{x}_1 . Since $x_{0 \rightarrow 2}$ are now known or imputed, we update the forward states $h_{0 \rightarrow 2}^f$ and predict $\hat{x}_3 = g(h_2^f, h_4^b)$. For
7 the second half, after predicting $\hat{x}_6 = g(h_4^f, h_8^b)$, we update $h_{4 \rightarrow 6}^f$ and $h_{8 \rightarrow 6}^b$. *Note that $h_{8 \rightarrow 6}^b$ have been updated before*
8 *when predicting \hat{x}_4 !* More generally, the forward states h^f are updated once whereas the backward states h^b are twice.
9 We encourage the reviewers to check the supplementary material, with code and visualizations of our decoding strategy.

10 **R2, R3 *Evaluation metrics.*** Evaluating generative models is an open problem, e.g., log-likelihood does not correlate
11 well with generation quality [Theis et al., 2015]. In our case, neither L2 nor log-likelihood can capture “realistic”
12 player behavior in basketball [Zheng et al., 2016, Zhan et al., 2019]. Hence, we follow previous work and compute
13 domain-specific metrics (speed, distance traveled, out-of-bounds rate) to compare trajectory quality. We will include
14 L2-loss for the basketball dataset, but note that NAOMI (0.013) still outperforms SingleRes (0.040).

15 **R1 *Motivation.*** In general, time-series data features different types of dynamics and missing value patterns compared
16 to text and images. Time-series data are often multi-resolution, which are exploited by our model via the divide-and-
17 conquer strategy. Note we do not use a fixed sampling scheme for missing values (see below). We would consider
18 combining NAOMI with convolutional or Transformer-based approaches to handle high-dimensional sequences.

19 ***Fixed length.*** No, our method does *not* assume fixed-length sequences. NAOMI can decode and train on varying-length
20 sequences, e.g., by padding shorter sequences to a maximal length.

21 ***Masking.*** We mask all dimensions for n randomly chosen time-steps, which is independent of the order of the divide-
22 and-conquer strategy (see Algo. 1). We used the same masking scheme for all methods, including MaskGAN and GRUI.
23 Note the “halving” scheme in Figure 2 is only an example: NAOMI is compatible with *any* masking pattern. If the
24 second half of a sequence is masked, NAOMI is pure forward inference (see supplemental material for results).

25 ***Auto-regressive baseline with divide-and-conquer.*** Note that “auto-regressive” and “divide-and-conquer” are mutually
26 exclusive decoding strategies, hence this baseline does not exist by definition.

27 ***Transformer.*** We agree that applying NAOMI to Transformer models is interesting, but leave this for future work.

28 **R2 *a bi-directed RNN ... to f_f and f_b encoders.*** NAOMI iteratively (re)-encodes and decodes as described above. Only
29 the initial sequence encoding (see Figure 2) behaves like a bi-directional RNN. ***Bi-cubic spline.*** We are happy to add
30 this, but since we compared with the state-of-the-art baselines (e.g., GRUI) for sequence imputation, we believe our
31 results stand on their own. ***Table 3, does Linear show smaller errors.*** The results in Table 3 are *not* errors, but *sample*
32 *statistics*, i.e., the closer to the Expert, the better. Linear has smaller values, but are actually worse as they are further
33 away from the Expert statistics. We chose these metrics for the reasons explained above.

34 **R3 *... might not address error propagation ...*** We agree that NAOMI may not fully solve error propagation for “any”
35 gap size between observed time steps. However, we compared many model variations and baselines on three time-series
36 datasets with different sequence lengths and varying missing-value proportions. Our extensive experiments have
37 empirically shown the effectiveness of NAOMI. We believe noisy coarse predictions are not an issue on these datasets
38 mostly due to the multiresolution structure and spatiotemporal smoothness of the data. Hence, even noisy coarse-level
39 predictions provide reasonable grounding points at finer resolutions.

40 ***Multiple decoders might be insufficient*** There appear to be several misunderstandings. Note that we randomly sample
41 masking patterns for each training step. Hence the model will see gaps of varying sizes during training. We compute
42 the loss at the sequence level to make sure the entire sequence look real, which co-trains the decoders. We repeatedly
43 use the reparameterization trick (L136) to make every sampling operation and hence the entire imputation procedure
44 differentiable. Our experimental results demonstrate this end-to-end training approach is sufficient for our datasets.

45 ***Maximum likelihood and adversarial training*** NAOMI is a non-autoregressive generator, *which can be trained with any*
46 *objective*. We used L2-loss for traffic and billiards because it is standard in those domains, see explanation above.

47 ***Missing values ... datasets.*** We will add analyses similar to Figure 6 for all datasets to the Appendix.

48 ***Minor comments.*** We will make Figure 4 more readable. For inputs of the backward encoder, the first dimension is 1
49 for known or predicted steps, 0 for masked steps, and the other dimensions are 0 for masked steps.

50 [1] E. Zhan, et al. (2019), “Generating Multi-Agent Trajectories using Programmatic Weak Supervision,” *ICLR*, 2019

51 [2] S. Zheng, et al. (2016), “Generating long-term trajectories using deep hierarchical networks,” *NIPS*, 2016