

1 We would like to thank the reviewers for the insightful remarks and comments. We address several concerns of the
 2 reviewers as follows.

3 **Main theoretical results and practical implications: R1 and R3.** Our main theoretical result is that the Max-Product
 4 algorithm can be exactly parameterized by the FGNN using number of parameters that are polynomial w.r.t. number of
 5 rank-1 tensors required to represent the higher order potentials (Corollary 7 in the main text). In worst case, it may
 6 require exponentially many parameters for arbitrary higher order potentials. However, we believe that there will be
 7 multiple practical applications that have reasonably sized representation. This will have to be verified empirically –
 8 we have demonstrated this for error correcting codes which showed improved performance under bursty channel, and
 9 for human motion prediction which showed state-of-the-art performance. Furthermore, in computer vision applications,
 10 commonly used handcrafted higher order potentials, *e.g.* Potts model or Robust Potts Model [DOI:10.1007/s11263-008-
 11 0202-0], can often be exactly or approximately represented using a moderate number of rank-1 tensors. FGNN is also
 12 advantageous when the potential functions are unknown and only the features corresponding to the factor are provided.
 13 In this case, we can use FGNN to learn the approximation of potential functions from data, where the size of the network
 14 can be decided based on the computational budget or amount of available data (to control over-fitting). As FGNN can
 15 exactly represent the Max-Product, it should be at least as powerful as the Max-Product algorithm. It can represent a set
 16 of algorithms including the Max-Product, so by learning from data it may potentially learn a better inference algorithm
 17 if Max-Product is not optimal. In our experiment, we show that FGNN outperforms the Max-Product algorithm.

18 **Possible over-fitting for synthetic data: R1.** In the synthetic data experiment, our goal is to train a better MAP
 19 inference algorithm, and thus the MAP assignment is used as the target for all neural network based algorithms. In the
 20 experiment, the number of all possible configurations is $2^{30} \approx 10^9$ for a length-30 binary chain structured MRF, and
 21 our training set only has 90000 items. In this case, it is unlikely that a neural network that simply over-fits the training
 22 set can have a good performance on the test set.

23 **Testing on novel graph structures: R3 and R1.** To address R3’s
 24 concern that the algorithm is tested only on the same graph struc-
 25 ture as seen in training, we conducted a new experiment to train
 26 the FGNN on fixed length-30 MRFs using the same protocol as
 27 Dataset3, and test the algorithm on 60000 random generated chain
 28 MRF whose length ranges from 15 to 45 (the potentials are generated using the same protocol as Dataset3). The result
 29 is in Table 1, which shows that the model trained on fixed size MRF can be generalized to MRF with different graph
 30 structures. This also further addresses the overfitting issue raised by R1.

Chain length	AD3	FGNN
(15, 25)	88.95%	94.31%
(25, 35)	88.18%	93.64%
(35, 45)	87.98%	91.50%

Table 1: Accuracy on dataset with different chain size.

31 **Factor without trivial representation of fixed dimension: R3.** In the experiment on human motion prediction, the
 32 potential functions are actually unknown and their representations are fully learned from data.

33 **Variational inference: R1.** The synthetic problems are MAP inference problems. As variation inference is designed
 34 for marginal inference, it cannot be directly used. Indirectly however, by using the zero-temperature technique, the
 35 variational inference can be applied to do MAP inference. By applying the zero-temperature technique, the objective
 36 function of the variational inference becomes equivalent to that of AD3 and MPLP. Particularly, the AD3 algorithm is
 37 guaranteed to converge to the optimal of that objective and thus it can be viewed as a variant of variational inference.

38 **Improve the writing: R1.** We will revise Section 3 to say that Corollary 7 is the main theoretical results of the paper
 39 and do careful proofreading. For Algorithm 1, Φ_{VF} and Θ_{VF} are parameters of the \mathcal{Q} and \mathcal{M} net in the Variable to
 40 Factor module, respectively. Φ_{FV} and Θ_{FV} are parameters of the \mathcal{Q} and \mathcal{M} net in the Factor to Variable module.

41 **Aggregation function: R2.** We choose the “max” aggregation function because theo-
 42 retically “max” is invariant to the duplication of a element in the set, while “sum” or
 43 “average” is not. In real applications such as human motion prediction, different factor
 44 may have different size, but for better parallelization we may need to pad all factor to
 45 the same size. In this case, we may simply duplicate a node in factor to do this. We replaced
 46 the “max” aggregation with “sum” aggregation in the LDPC experiment and typical
 47 result is shown in Figure 1, where both algorithm achieve almost the same performance.

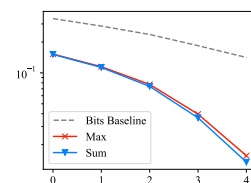


Figure 1: Comparison of “sum” and “max” aggregation.

48 **Training time: R2.** The training time is in Table 6 of the supplementary file.

49 **Representativeness of the LDPC data: R2.** In this experiment, as previous work [Kim et al., 2018, Zarkeshvari
 50 and Banihashemi, 2002] we do not have any assumption on the input signal, thus both the train and test set the input
 51 signal are uniformly sampled from $\{0, 1\}^{48}$. We can add error bars in the next version to indicate how well the results
 52 represent uniformly distributed inputs.

53 **Relation with Yoon et al. and publishing code: R3.** We will add discussion on relation with Yoon et al., which is on
 54 pairwise potentials, as well as detailed experiment settings and the url to the source code in the next version.