

1 We thank all reviewers for their time and thoughtful comments. **Summary of positive reviews:** Reviewers liked our
2 proposed algorithm (BPNN) because it is innovative (@R1), elegant (@R3), interesting (@R4), and theoretically
3 grounded (@R2, @R4) with clear proofs (@R3). Our approach shows solid experimental improvements for computing
4 the partition function of Ising models, community detection problems, and factor graphs representing boolean formulae
5 (@R2, @R3, @R1). Additionally, we demonstrate robustness to perturbations in the test distribution (@R1). Reviewers
6 thought that our paper was generally well written (@R2) and that it does a good job of explaining our contributions
7 and their significance (@R3), although it is inaccessible to readers who aren't familiar with the variational view of
8 BP (@R3). **Summary of negative reviews:** The most serious concerns were that it was challenging to understand our
9 model's structure/output (@R4) and the training protocol we used (@R1). Additionally, there were questions/confusions
10 about our experimental setups.

11 **@R1, @R4: What is the output of BPNN?** BPNNs consist of two parts. First, iterative BPNN layers output messages,
12 analogous to standard BP. These messages are used to compute beliefs using the same equations as for BP. Second, the
13 beliefs are passed into a Bethe free energy layer (BPNN-B) that outputs an estimate of the log partition function. This
14 layer generalizes the Bethe approximation by performing regression from beliefs to $\ln(Z)$. Alternatively, when the
15 standard Bethe approximation is used in place of the (trainable) BPNN-B layer, BPNN provides additional guarantees.

16 **@R1: What is the training protocol?** In all our experiments, we initialized the BPNN to output the Bethe approxi-
17 mation obtained by running BP for a fixed number of iterations. We used the mean squared error between the BPNN
18 prediction and the ground truth log partition function as our training loss. When training BPNN-D without a BPNN-B
19 layer, we ran BPNN-D for a random number of iterations between 5 and 30 at every training step to encourage quick
20 convergence to a fixed point. (When training with a BPNN-B layer, iterative layers were applied for a constant number
21 of iterations, as convergence to a fixed point is no longer required.) We will highlight these details in the final copy.

22 **@R4: When is convergence to a fixed point unnecessary? How was the BPNN-B layer designed?** We can think
23 of BPNN as a trainable computation graph (neural network) that mimics the standard (unrolled) BP computation.
24 Convergence to a fixed point is unnecessary for good predictive accuracy; however, this extra flexibility makes theoretical
25 analysis more difficult. The Bethe layer (BPNN-B) was carefully designed to generalize the Bethe approximation while
26 maintaining invariance to factor graph isomorphism (please see Lemma 4 in the appendix).

27 **@R1, @R4: What version of BP is compared against? What does "standard BP" refer to?** We compared against
28 BP where we tuned the traditional damping coefficient and message update strategy (sequential or parallel) for best
29 results. "Standard BP" refers to belief propagation run with traditional damping but without learned modifications to
30 messages (via the operator $H(\cdot)$).

31 **@R4: Did you compare with a strong, baseline algorithm that utilized the same training dataset as BPNN?** Yes,
32 we compared with and found significant improvements over the Graph Isomorphism Network GNN.

33 **@R3, @R4: Did your experiments include problems where BP converges well?** Yes, BP converged on all commu-
34 nity detection problems and 88% of Ising models (line 178).

35 **@R1: "Why would we expect BPNN-D to converge to better local optima than BP? Because the learned H
36 function helps avoid local optima? If this is true, it is an interesting observation by itself."** Yes, that is the intuition
37 and we agree it is interesting! Fixed points of BP are local optima of the Bethe free energy. By training BPNN-D to
38 predict the exact partition function, BPNN-D can learn to find better local optimum.

39 **@R1: "I'm very curious how overfit (if at all) these models are to the domain on which they are trained."** We
40 explored cross domain generalization in our propositional model counting experiments and found that the quality of
41 results is dependent on the similarity of domains. We found that we could learn improvements that translate across a
42 broad variety of domains, although the improvements were less dramatic (lines 289-293).

43 **@R4: "The theorems in the paper ... do not add anything in understanding when BP on loopy graphs converge"**
44 That is correct, our theorems make no claims about the convergence properties of BP on loopy graphs. Our theorems
45 are given to *precisely characterize the relationship between fixed points of BP and BPNN-D* in terms of properties of
46 the learned operator $H(\cdot)$.

47 **@R4: "The theorems in the paper are essentially trivial ... Without going through all provided proofs in detail,
48 the theoretical justification of the stated theorems seems to be correct."** Even though the theoretical results may
49 seem intuitively correct in hindsight, they provide valuable insights into BPNN. We think it is unfair to characterize the
50 theorems as trivial because similar results do not hold for standard GNNs.

51 **@R3, @R4: Can BPNN be used to estimate marginals?** Yes. We focused on estimating the partition function
52 because these estimates are better understood for BP. However, in a small set of experiments on estimating marginals for
53 community detection, we found large improvements over BP and GNNs. This is a promising direction for future work.

54 **@R3: Can the BPNN-B layer be modified for computational efficiency on high dimensional factors?** Yes. The
55 BPNN-B layer can be modified to use an attention mechanism for efficient invariance given high dimensional factors.

56 **@R2: Will you open source your code after publication?** Yes!