

1 We thank the reviewers for their insightful feedback; we address each review below.

2 **R1: “...it is known how to solve the BFE optimisation problem by double loop algorithms”** Our understanding is  
3 that double-loop algorithms (Yuille, 2001) converge to a local minimum, but not necessarily a global one. Further,  
4 double-loop algorithms outperform LBP generally when LBP does not converge, but in our experiments LBP converged  
5 reliably, and we thus anticipate that double-loop optimization will not necessarily improve further. We also note that  
6 despite the superiority of double loop algorithms, LBP remains popular, and so amortizing it seems useful in our view.

7 **“...what is meant by ‘they are run once...’”** A Transformer or RNN is run over a sequence of embeddings corre-  
8 sponding to the nodes in the graph to obtain a sequence of representations. We do *not* then further update these  
9 Transformer/RNN representations with message-passing style updates as in much recent work in graph neural networks.

10 **“...meaningless for pairwise marginals...”** Agreed. We included the pairwise marginals just for completeness.

11 **“Ising model...expect that the estimation quality will degrade with the (average) interaction strength.”** Thank  
12 you for this suggestion! We tested our approach in settings where the interaction strength is higher by sampling pairwise  
13 potentials from  $\mathcal{N}(0, 3)$  (unary potentials are still sampled from  $\mathcal{N}(0, 1)$ ). For a 10x10 grid we obtain correlations  
14 of 0.460, 0.641, and 0.770 for Mean Field, Loopy BP, and the inference network, respectively. For a 15x15 grid we  
15 analogously obtain 0.435, 0.665, and 0.738. When sampling from  $\mathcal{N}(0, 5)$  we obtain 0.358, 0.525, 0.619 for these 3  
16 approaches in the 10x10 case, and 0.362, 0.459, 0.564 in the 15x15 case.

17 **“The experiments have in my view a preliminary character”** We agree our experiments are on small datasets. Our  
18 contribution is mostly methodological and a first step toward learning undirected models in this way, in line with early  
19 work for learning directed models with VAEs. We also note that we consider 3 fairly different archetypal settings.

20 **“some details of them are missing”** We apologize for the missing details, which had to be relegated to the supplemen-  
21 tary materials due to space. We will make sure to include them should the paper be accepted.

22 **R2: “...having some higher order factors”...** Having too many higher order factors could be slow (since marginals  
23 scale exponentially in the order) but having a few may improve performance; we will investigate this!

24 **“...optimizing in the subspace...”** We explored this approach in preliminary experiments, and while it can succeed it is  
25 much less scalable: it requires (pre)computing the SVD of a large matrix (i.e.,  $O(|\mathcal{F}|^2)$ ) and we need an SVD for each  
26 graph topology (c.f., HMMs, which have different lengths per sentence). We will add more discussion around this point.

27 **“...hardware...”** We ran our experiments on 1080 Tis, using pytorch 1.1.

28 **“Why are the inference networks different...”** Empirically we observed small differences between architectures for  
29 different tasks in preliminary experiments, and we chose the best; see answer to R3 as well.

30 **R3: “The marginal constraints can also be enforced by matrix operations..”** See response to R2.

31 **“The numbers in Table 1 and Figure 2 seem inconsistent...”** The figures in Table 1 are for both unary and pairwise  
32 marginals. While LBP does slightly better for unary marginals, the inference network does much better for pairwise,  
33 and when averaged across both in Table 1, the inference network does better overall; we will clarify this further.

34 **“...I wonder whether some architecture of inference networks are better than others...”** We experimented with  
35 both Transformers and RNNs, and picked the architecture with the best performance in preliminary experiments. For  
36 HMMs, we find Transformers give an NLL of  $\approx 116$ , slightly worse than RNNs; we will include these results.

37 **“The form of discrepancy used ... is not clear.”** We apologize for the lack of clarity: we used L2 for the synthetic  
38 experiments and Jensen-Shannon divergence for the rest, though L2 was only slightly worse. We will add this detail.

39 **“...random seed as a hyperparameter...”** In our experience discrete latent variable models can be quite sensitive to  
40 initialization, and so to fairly compare all model variants we randomly sample the same number of seeds for each; we  
41 do not regard this as a hyperparameter but as giving each model/method equal opportunity in experiments.

42 **“...did not get numbers similar...”** The following reproduces our numbers on a 1080 Ti, using pytorch 1.1, by epoch  
43 9: `python pen_uhmm.py -cuda -K 30 -bsz 32 -dropout 0.3 -ilr 0.0003 -infarch rnnnode -init 0.001 -just_diff -lamb_size 64 -loss`  
44 `alt3 -lr 0.0001 -markov_order 3 -max_len 30 -not_inf_residual -optalg adam -penfunc js -q_layers 1 -qemb_size 150 -qinit 0.001`  
45 `-seed 21442 -vemb_size 64 -epochs 10`. The code prints out *perplexity* rather than the NLL reported in the paper, and thus  
46 looks higher. We apologize for the potential confusion. (To convert, take  $\log(\text{perplexity}) \times 50509/2747$ ).

47 **“...the authors probably ran over 100k seeds...”** We did not run 100k seeds; all methods were given 100 random  
48 configurations, as the code (line 609 in `pen_uhmm.py`) makes clear. Regarding stability, we calculated the std. deviation  
49 of best NLL obtained per run for each method in Table 4, giving (from top to bottom): 5.3, 0.7, 1.1, 0.5, 6.4, 2.8, 1.1.  
50 Thus amortized BFE is more stable than exact inference and LBP, but less than Mean-Field+BL and 1st Order.