We thank the reviewers for their comments.

**Reviewer #1**: Regarding the challenges that arise in proving the structural result: Compared to structural results on learning MRFs, such as [2], the main difference is that we obtain a time complexity for learning RBMs that is believed to be impossible for learning MRFs, by exploiting the latent-variable structure of RBMs. The main challenge is to relate the properties of the induced MRF to this latent-variable structure. The task is non-trvial due to the non-linear effect that latent variables have on observed variables.

To overcome the challenge, the first step is Lemma 6, which is an "interchange of sums" argument that expresses the mutual information on the observed variables as a sum over the latent variables. This expression is convenient in its explicit relation between observed variables and latent variables, but it is still complicated. Then, the second step is Lemma 7, which limits the cancellations that can happen in a sum of products by the number of terms in the sum. Combined, these two steps allow us to constrain the cancellations of the low-order interactions between the observed variables, in terms of the latent-variable structure of the RBM.

Regarding conditioning on smaller subsets of nodes in Theorem 11: Such bounds are plausible, because the conditional probabilities with smaller subsets are simply weighted averages of the full conditional probabilities. However, there is an error both in the terms and in the weights of these weighted averages, and we do not currently have bounds in this direction.

**Reviewer #2**: Regarding the phrase "few latent variables": We agree that the phrase is not sufficiently accurate. Perhaps a better option would be "locally sparse latent variables". We intend to switch to this phrase, or some similar one, in the final version of the paper. We will also ensure that the assumption is made precise earlier in the exposition.

Regarding the properties leveraged to obtain the improved bounds: Some intuition is as follows. The difficulty of estimating the Markov blanket of an observed variable lies in the fact that the low-order interactions between the observed variables can vanish. If that were not the case, a greedy $\tilde{O}(n^2)$ algorithm would work (this is indeed the case if all the RBM weights are positive [1]). Then, our analysis shows that the order up to which all interactions can vanish is constrained by the number of latent variables. Or, in other words, that distributions in which all interactions up to a large order vanish are "complex" and require "many" latent variables.

Regarding other approaches that would estimate the Markov blanket: Unfortunately, estimating the Markov blanket of an observed variable seems to be exactly where the difficulty lies. Based on all prior literature, estimating the Markov blanket would only be guaranteed to scale as $\tilde{O}(n^d)$. Then, our guarantee improves this to $\tilde{O}(n^{2^s+1})$.

Regarding the result in Section 5: In case it was not sufficiently emphasized in the paper, we note that the sample complexity, and hence also the time complexity, of structure learning for RBMs necessarily depends on the minimum potential, due to information-theoretic arguments [3]. All prior algorithms have this dependence. Then, the contribution of Section 5 is that we show a way to eliminate this dependence; we necessarily lose the guarantee on structure recovery, but we still guarantee accurate conditional prediction.

Regarding the definition of $U$ on line 203: The sentence as written is correct. It turns out that it is possible to "factor out" (into $\bar{f}$) the contribution of latent variables that are connected to $u$ but not to any observed variables in $I$. Note that our definition of $s$ also excludes the latent variables that are connected only to the observed variable in question.

**Reviewer #3**: Regarding relevance: Sparsity, in particular as a bound on the maximum degree of a node, is a common assumption in the analysis of graphical models. Note that, by itself, a bound of $d$ on the maximum degree in an RBM implies that $s \leq d^3$. Then, our assumption can be seen as simply a strengthening of this particular aspect of sparsity, such that few latent variables act on the neighborhood of each observed variable.

More generally, our main motivation behind this work was to carve out good sets of assumptions for which efficient learning is possible with time complexity less than $\tilde{O}(n^d)$. This is the first work that achieves such a time complexity without ferromagneticity, i.e., with arbitrary weights in the RBM.

# References

[1] Guy Bresler, Frederic Koehler, and Ankur Moitra. Learning restricted Boltzmann machines via influence maximization. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 828–839. ACM, 2019.

[2] Linus Hamilton, Frederic Koehler, and Ankur Moitra. Information theoretic properties of Markov random fields, and their algorithmic applications. In *Advances in Neural Information Processing Systems*, pages 2463–2472, 2017.

[3] Narayana P Santhanam and Martin J Wainwright. Information-theoretic limits of selecting binary graphical models in high dimensions. *IEEE Transactions on Information Theory*, 58(7):4117–4134, 2012.