

1 We thank the reviewers for their insightful comments and interest in the work, find below line by line responses. For  
2 space reasons we do not reply to comments on organization or typos but have incorporated them into the paper.

3 **R1, R2, R3, R5: Limitation to trees.** Several reviewers raised questions about our limitation to tree graphs. We  
4 apologize for not motivating our assumption more clearly, and will modify the paper to clarify that the algorithm works  
5 whenever the undirected components of the graph form a forest. As briefly mentioned in L51-52, after the observational  
6 stage, the undirected components of the essential graph are chordal. Orienting each of them provides no information on  
7 the others, thus they have to be handled separately by any algorithm. Our algorithm only requires these components  
8 be trees, but they do not need to be connected. Note we can identify all v-structures from observational data. For  
9 example this assumption is satisfied when the original graph is bi-partite, since chordal components of bi-partite graphs  
10 are forests. Examples of bi-partite causal graphs occur in system biology networks, e.g. gene-disease networks or  
11 gene-protein networks [<https://www.ncbi.nlm.nih.gov/pubmed/27265032>]. We will add this discussion to the paper.

12 **R2: ...information greedy algorithm is suboptimal. (low significance)** While we don't dispute this, we point out  
13 that the dramatic failure of information greedy puts our problem in a different class from many active learning scenarios.

14 **[the case in which]  $P(Y = y)$  is very similar to  $P(Y = y|X_i = 1)$ .** As R2 correctly states, our algorithm may not  
15 be optimal if some edges are much weaker than others (i.e. when Condition 1 in Sec. 3.4 does not hold). Note that  
16 under Condition 1, we provide formal guarantees for a version of our algorithm. We also believe that Condition 1 is  
17 justified in real-world examples: this is because with finite observational data, we expect weak edges to be left out by  
18 methods that learn the graph skeleton.

19 **...how such an information greedy algorithm works w.r.t. finite-sample case.** In the finite sample case, the  
20 information greedy algorithm must take more interventions than in the infinite sample case. Hence the number of  
21 required interventions is still linear in the size of the graph  $n$  for our counter example structure. Meanwhile under  
22 Condition 1, the infinite-sample central node algorithm can be applied to the finite sample case by repeating each  
23 intervention until a branch has at least  $1 - \delta/(\log n)$  probability, which requires order  $\log \log n$  repetitions in expectation  
24 by Prop. 2. By the union bound then, this algorithm finds the solution with  $1 - \delta$  confidence in number of interventions  
25 nearly logarithmic in  $n$ , a major speedup over information greedy.

26 **...incorporate a prior distribution...** We only consider priors over the graph structure. Our approximation guarantees  
27 with respect to the optimal algorithm hold for any prior assumed over the source node locations (see Sec. 3.1).

28 **R5: interventions are assumed to be single target. Is there any straightforward way to extend the updates...to  
29 simulations interventions?** Lemma 1 can indeed be extended, we did not originally include it because it would  
30 complicate the paper and raise several corner cases which complicate the algorithms. We will add it to the supplement.

31 **The type of interventions used in this work is Pearl's atomic intervention...it may be possible that the cause  
32 affects the effect only under certain outcomes.** We wrote the paper with  $do(X_i = 1)$  for notational simplicity, as  
33 long as there is any value  $a_i$  for which  $do(X_i = a_i)$  affects the effect variables, we can find it from the observational  
34 data and the theory will still hold. We will expand and rigorize this discussion in the camera-ready.

35 **the extension is not clear...In the experiments, the authors have restricted themselves to binary variables.** We  
36 will add a discussion and new experiments for the non-binary variables, the results are qualitatively the same. Essentially,  
37 the extension hinges on the fact that we simply need the result of the do-intervention to provide information on the  
38 direction of the root via the posterior update (Lemma 1), all our results follow from this.

39 **Connection with [Ghassami et al. 2017]:** Thank you for pointing us to this paper, we will be sure to cite it and make  
40 the connection clear. Indeed, their Definition 8 and our central node definition are the same, and our algorithm can be  
41 seen as repeating the first step of ProBAL with evolving posteriors. In fact, this is true for any non-adaptive algorithm  
42 whose first step is intervening on the central node of the prior. However, ProBAL is a non-adaptive algorithm designed  
43 for a fixed budget setting (number of interventions). Interestingly, our results show that it is a good algorithm with  
44 a constant factor approximation guarantee even in the adaptive, fixed confidence noiseless setting. We also show an  
45 adaptation of this works well in case of restricted nodes and in the noisy case with an assumed bound on the noise level.

46 **Comparison to [Hauser and Buhlmann, 2014]:** While OptSingle is non-adaptive, it can be turned into an adaptive  
47 algorithm if one applies it repeatedly on the remaining essential graph after the intervention and application of Meek  
48 rules. The objective is to minimize the number of unlearned edges in the worst case graph. In the case of a tree with  
49 a uniform prior, this algorithm will also intervene on the central node, since only that will minimize the worst case  
50 number of unoriented edges. However, when the prior is non-uniform the two algorithms are not equivalent: OptSingle  
51 will only use the structural properties and not take into account the prior. In the noisy case, a single intervention will not  
52 change the support of the posterior distribution, hence OptSingle will continue intervening on the same node without  
53 making any progress. Hence we did not see a straightforward way to compare to it in our experiments, all of which  
54 consider the noisy case. We will add appropriate citations and discussion of these papers in the camera-ready.