1 Thank you to all reviewers for the very careful feedback. We respond here, and **we will update the paper accordingly**!

2 (**Experiments**) We evaluated additional models, providing further empirical evidence on the PGNs' benefits: (**SupGNN**)
3 baseline GNN whose latents have been supervised to be predictable of masks and pointers (without explicitly using this
4 information); (**PGN-Asym**) PGN without symmetrising pointers; (**PGN-MO**) PGN where only masks are supervised
5 (and pointers are not directly optimised); and (**DGM**) [19], a recent strong baseline method on latent graph inference.
6 Here we provide results on the **hardest** test set (LCT; $n = 100$; ops $= 150$); results on other sets mirror these findings.

| | GNN | SupGNN | DGM | PGN-MO | PGN-Asym | PGN |
|---|---|---|---|---|---|---|
| 7 | $0.401_{\pm.123}$ | $0.541_{\pm.059}$ | $0.524_{\pm.104}$ | $0.558_{\pm.022}$ | $0.561_{\pm.035}$ | $\mathbf{0.616}_{\pm.009}$ |

8 (**R1 / R2 / R3**) Our work evaluates the extent to which existing neural networks are potent reasoning systems, and the
9 minimal ways (e.g. inductive biases / training regimes) to strengthen their reasoning capability. Hence our aim is not to
10 outperform classical algorithms, but make their concepts accessible to neural networks. PGNs enable reasoning over
11 edges not provided in the input, simplifying execution of any algorithm requiring a pointer-based data structure. PGNs
12 can find direct practical usage if, e.g., they are pre-trained on known algorithms and then deployed on tasks which may
13 require similar kinds of reasoning (with encoders/decoders "casting" the new problem into the PGN's latent space).

14 (**R1**) Inspired by your review, we scaled up the LCT test set to $(n = 200, \text{ops} = 300)$ where the PGN $(0.636_{\pm.009})$
15 catches up to Oracle-Ptr $(0.619_{\pm.043})$. We hope that this illustrates not only the robustness of PGNs to larger test sets,
16 but also provides a quantitative substantiation to our claim about ground-truth LCTs not having favourable diameters.

17 (**R1**) Our contribution is *both* in the inductive biases (e.g. pointers/masks) proposed and the provided supervision.
18 To strengthen this claim, we compare against **SupGNN** in experiments above. Supervising provides a clear boost in
19 performance, but is still insufficient to outperform PGNs, showing the importance of the architectural choices as well.

20 (**R2 / R1**) To strengthen our empirical contribution on learning latent graphs, we also compare against the recently
21 proposed **DGM** [19] loss function in the experiments above, finding PGNs extrapolate much better to larger test sets.

22 (**R2 / R4**) We added experiments on asymmetric pointers (**PGN-Asym**), with a clear drop in performance on the LCT
23 task. Corresponding to R4's comment, we find that while symmetrising the pointers is not strictly necessary, empirically
24 it is helpful because it makes the inferred graph slightly less sparse, enabling us to rectify for the case where the pointer
25 mechanism makes mistakes without sacrificing sparsity. In this sense, a more "global" view of each node is desirable.

26 (**R3**) Our work extends [42] by the addition of the **pointer / masking** inductive biases, allowing one to more efficiently
27 simulate algorithms for which the input graph does not correspond to the actual links used for reasoning, or cases where
28 an input graph is not provided at all. As we acknowledge, while conditional masking had been introduced in [50], here
29 we utilise masking to update node state rather than exclude it from the input—which can reduce information loss.

30 (**R3**) We now compare against **PGN-MO** in the experiments above, which supervises on **m**asks **o**nly and does not
31 directly optimise pointers. The results show that we can relax most of the supervision constraints and still outperform
32 baseline approaches (such as GNN and DGM), while for full returns an explicit pointer-based loss is required.

33 (**R3**) Please note that PGN-Ptrs requires *two passes*: first training a PGN, then re-using and fixing its pointers on a
34 second training pass. Allowing for dynamic adjusting of pointers works better on complex algorithms with less margin
35 for error, such as the LCT setup. If the PGN mechanism makes a mistake at any epoch, PGN training can meaningfully
36 recover from this while PGN-Ptrs don't have this flexibility. For DSU, path compression yields very shallow data
37 structures that are modellable in many equivalent ways, hence such errors are less punishing for PGN-Ptrs.

38 (**R3**) Thank you for remarking the correctness of our "expressivity" claim—we will amend as you suggested!

39 (**R3**) Various answers: $F_1$ measures performance across the (binary) query responses $y^{(t)}$, handling class imbalance. As
40 mistakes are measured on $y^{(t)}$, errors at step $t$ do not imply errors at step $t + 1$. By parallelisable, we imply that they
41 allow for using GPUs more effectively (e.g. the rollout vs. ground-truth in Fig. 4). 1-NN refers to 1-nearest neighbours,
42 for which efficient algorithms exist and we don't need to compute all $O(n^2)$ $\alpha_{ij}$ values. $\|$ is concatenation.

43 (**R4**) We agree that multiple GNN steps may be theoretically required to reconcile our work with teacher forcing the
44 data structure. We found it sufficient to consider only one-step here, but as tasks get more complex—especially when
45 compression is no longer applicable—dynamic number of steps (e.g. function of dataset size) is likely to be appropriate.

46 (**R4**) Thank you for suggesting the GAT experiment! While it seemed to compare favourably to GNN/DeepSets on the
47 (more compressible) DSU task, the unrestricted GAT model failed to get off the ground at all for the LCT task.

48 (**R4**) Thank you for the minor comments, which we will incorporate! Especially for pointing out the work of Goel *et*
49 *al.*, which allows us to report a tighter analysis on time complexity (as we effectively use randomised linking-by-index).