

1 We thank all three reviewers for their very thoughtful and detailed comments, with which we largely agree. Some
2 requested additions (e.g., details about experiments and interventional datasets, an example of an empirical evaluation)
3 were present in previous versions of the paper and were removed for space. We will add these sections to any revised
4 version, either in the main body of the paper or the Supplemental Material. More details appear below.

5 **Additional citations (R1)** We thank Reviewer 1 for the highly relevant pointers to the statistics literature—these will
6 be cited and discussed in any revised version. Within-study comparisons provide a great source of data that is
7 both experimental and observational and should be discussed as a source of interventional data.

8 **Originality (R1)** Some of the points discussed in this paper have been touched on in previous work in statistics and
9 quantitative social science (e.g., the citations mentioned by Reviewer 1). However, this work is not widely
10 known in the machine learning community (the primary audience for NeurIPS), and this work often concerns
11 a highly specific evaluation context (i.e., single-treatment/single-outcome, average treatment-effect, etc.). As
12 seen by our survey, only a small fraction of the ML community performs evaluations that are both empirical
13 and interventional, and there are very few standard data sets in the field that allow for this. We believe that
14 further evidence, targeted at researchers within computer science and addressing the approaches and concerns
15 of that community, is necessary to promote the wider adoption of more principled evaluation methods.

16 **Evaluation guidance (R1)** Reviewer 1 notes that the paper provides little guidance on how to perform a principled
17 empirical evaluation. An example of an empirical evaluation was present in a previous version of the paper,
18 but it was cut for space. It is clear that, for the type of evaluation we advocate to be better understood and
19 adopted, more description and guidance is needed. We will include a detailed example in any revised version.

20 **Interventional vs Structural measures on the correct structure (R2)** Reviewer 2 is correct that, for a correctly
21 learned structure, as long as the parameterization is correct, interventional and structural measures should
22 agree. Interventional measures are valuable when the learned structure is only approximately correct, and
23 empirical evidence indicates that this is almost always the case. We will add more discussion of this to any
24 revised version to make this clearer. Some related experiments were moved for space to the Supplemental
25 Material, showing how structural Hamming distance and structural intervention distance penalize over- and
26 under-specification differently than total variation distance.

27 **Untested influences (R2)** The critique of untested influences refers to potential "unknown unknowns" in the data-
28 generating system. In many real-world systems, even with strong domain knowledge, there often exist factors
29 that are outside the researcher's knowledge. This is generally not possible in a synthetic system. While latent
30 variables can be added, they are still defined and created by the researcher, limiting the realism of the data. We
31 will clarify this in any revised version.

32 **Number of samples in synthetic data experiments (R2)** The number of samples was chosen to match the number of
33 samples available from the software system. This is 2599 for networking, 473 for jdk, and 5000 for postgres.
34 We acknowledge there may be issues with low sample sizes and will discuss this in any revised version.

35 **Structural measures implicitly assume a DAG (R2)** Our intent with this statement was that some structural measures
36 can only be used by algorithms that produce a DAG as output. There are many ways this assumption could
37 be violated. Acyclicity is one such violation, which could hinder the use of certain structural measures in
38 evaluation of an algorithm that outputs a cyclic graph. Another possible violation is an algorithm that does not
39 output a graphical model, such as a recent work on learning probabilistic programs. We will clarify this intent
40 in any future version.

41 **Total variation distance calculation (R2)** In Figure 2, TVD is calculated as the distance between the learned distribu-
42 tion and the true distribution, just as SHD is calculated as the distance between the learned structure and the
43 true structure. Here, the true distribution is known because the data is synthetically generated, although such
44 true distributions can also be obtained from interventional experiments as they are elsewhere in the paper.

45 **Domain-specific simulations (R3)** Domain-specific simulation systems are a very useful approach to consider, and
46 we will add a discussion of them in any future version of the paper. We already mention one such system (the
47 DREAM in-silico challenges), but we will add a far more general discussion of this approach. A sufficiently
48 sophisticated simulation falls on a spectrum between purely synthetic and purely empirical data, and such
49 simulations are valuable because they are often highly complex, are created by someone other than the
50 researcher, and are created for a purpose other than evaluation.

51 **Capabilities and limitations of current interventional datasets (R3)** We agree that a more detailed discussion of
52 the currently available interventional datasets would be highly beneficial. This sort of discussion was omitted
53 for space but will be included in at least the Supplemental Material of any revised version.