

ADDED AND MODIFIED CONTENT FOR PAPER #575

Anonymous authors

Paper under double-blind review

1 ADDED NEW DEFINITION

2 **Definition 100** (Causal representation). *The causal representation \hat{X} represents the computed values*
 3 *of causal variables when constructing a causal model, i.e., the quantified values from causal variables.*
 4 *Causal representations should meet the following two conditions:*

- 5 • *Correlation Condition: For any two causal variables that are not independent, their corre-*
 6 *sponding causal representations must be correlated.*
- 7 • *Causation Condition: For any two causal variables that have direct causal relationship,*
 8 *their causal representations should contain not only the information about their correlation*
 9 *but also information about the causal relationship.*

10 *For example, in SCM, information about the noise terms can be included, where correlation rela-*
 11 *tionship can be determined by fitting or comparing measures (such as cosine similarity), and causal*
 12 *relationship can be determined by examining the residuals Σ of the fitted model (Chen et al., 2023a).*

- 13 • $\Sigma_X \perp\!\!\!\perp Y, \Sigma_Y \not\perp\!\!\!\perp X \Rightarrow Y \rightarrow X$
- 14 • $\Sigma_X \not\perp\!\!\!\perp Y, \Sigma_Y \perp\!\!\!\perp X \Rightarrow X \rightarrow Y$
- 15 • $\Sigma_X \not\perp\!\!\!\perp Y, \Sigma_Y \not\perp\!\!\!\perp X \Rightarrow L \rightarrow X, L \rightarrow Y$
- 16 • $\Sigma_X \perp\!\!\!\perp Y, \Sigma_Y \perp\!\!\!\perp X \Rightarrow X \rightarrow L, Y \rightarrow L$

17 The difference between causal representation and ordinary deep representation lies in the “Causation
 18 Condition” mentioned in Definition 100. In general, deep representation can only meet the Correlation
 19 Condition, meaning it can only identify correlations. However, causal representation can identify not
 20 only correlations, but also causal relationships.

21 We noted that some reviewers also had questions about the causal structure. By causal structure, we
 22 mean causal diagrams, which we believe is a widely recognized term.

23 2 THE DETAILS ABOUT CAUSAL CONSISTENCY

24 Let’s illustrate the concept of causal inconsistency using a simple example: Consider a simple causal
 25 structure $A \leftarrow C \rightarrow B$, from the structure, we can see that A and B are correlated. This is due to
 26 $P(A, B|C) = P(A|C) * P(B|C) \Rightarrow A \not\perp\!\!\!\perp B (A \perp\!\!\!\perp B|C)$. In the case of single-valued variables, we
 27 can conduct independence tests on all samples of A and B to ascertain whether they are correlated. If
 28 they are, we then conclude that the causal structure and representation are consistent; otherwise, they
 29 are inconsistent. For multi-valued variables (such as deep representation satisfying Definition 100),
 30 one method to approximate this “correlation” is using cosine similarity or mean squared error (MSE).
 31 If the representations of A and B are similar, we also consider the structure and the representation to
 32 be consistent.

33 Hence, we use a similarity matrix to measure this “inconsistency.” To continue with the example
 34 above, we hypothesize two similarity matrices for the structure and representation respectively,
 35 $Sim^s \in \mathbb{R}^{3 \times 3}$ and $Sim^r \in \mathbb{R}^{3 \times 3}$. In these, $Sim^s_{i,j} = P(i|j)$, and $Sim^r_{i,j} = \text{cossim}(i, j)$. The MSE
 36 between these two similarity matrices Sim^s and Sim^r is used to measure inconsistency - if the MSE
 37 is close to 0, it indicates that Sim^s is approximately equal to Sim^r , i.e., the causal structure and the
 38 causal representation are essentially consistent, otherwise they are inconsistent.

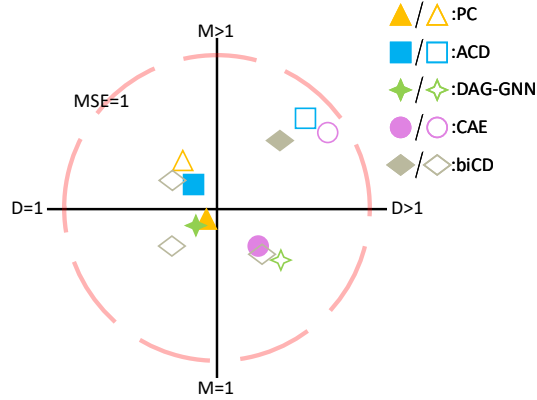


Figure 100: We compared the consistency of different methods in 3 data paradigms (if available). The consistency was represented by the MSE of the similarity matrices for structure and representation. The filled markers represent methods being in their default data forms, while the hollow markers signify that they are in extendable but non-default data forms.

39 3 MODIFIED SECTION: WHY CAUSAL INCONSISTENCY ARISES?

40 Figure 100 visualizes evaluation results of causal consistency via tested 5 methods: PC (Kalisch &
 41 Bühlman, 2007), ACD (Löwe et al., 2022), DAG-GNN (Yu et al., 2019), CAE (Chen et al., 2023a),
 42 and biCD (Chen et al., 2023b). They represent prevalent methods in specific data forms, respectively.
 43 Two conclusions can be obtained from Figure 100:

- 44 • The strongest causal inconsistency is found in indefinite data forms ($M > 1, D > 1$), while
 45 definite data ($M = 1, D = 1$) performs the weakest causal inconsistency.
- 46 • When existing methods are applied to non-default data forms (hollow markers), their
 47 consistency performance is always inferior to the native methods for that data form.

48 In addition to experimental results, we also provide comprehensive theoretical analysis for three
 49 different data paradigms:

- 50 • **Definite Data ($M=1$ and $D=1$):** The causal strength f is fixed and can be recovered through
 51 statistical properties in the data (for example, independent tests, independent component
 52 analysis, rank of covariance, etc.). Therefore, the estimated causal representation $\hat{X} = X$
 53 ($D=1$), and the causal strength $\hat{f} = f$ ($M=1$). The subtle inconsistencies in Figure 100 arise
 54 from biases or confounding in the sampling process.
- 55 • **Semi-Definite Data ($D>1$ and $M=1$):** According to Definition 100, there are differences
 56 between the causal representation \hat{X} and the input representation X , therefore we can't
 57 directly optimize causal representation through $loss(\hat{X}, X)$. Fortunately, in this situation, f
 58 is fixed, so we can map \hat{X} to a unique \hat{f} without parameters, and then optimize the process of
 59 causal discovery through $loss(\hat{f}, f)$. Thus, \hat{f} is the projection of \hat{X} and inherently consistent.
 60 The effectiveness of \hat{X} comes from: $\hat{X} \Leftrightarrow \hat{f} = f \Leftrightarrow X$. The slight inconsistencies in
 61 Figure 100 result from biases of projection.
- 62 • **Semi-Definite Data ($M>1$ and $D=1$):** The estimated causal strength \hat{f} can be viewed as
 63 distribution determined by X and encoder parameter φ , denoted as $\hat{f} = h(X, \varphi)$, and is
 64 optimized through $loss(\hat{f}, f)$. \hat{X} can be estimated via inverse function h^{-1} , because when
 65 the causal variable does not need to be quantified into a deep representation, there exists
 66 an error $loss(\hat{X}, X)$ such that $\hat{X} = X$. From this, we can get the equivalent equation:
 67 $\hat{f} = f \Leftrightarrow X = \hat{X}$ (\Leftrightarrow is because of the ground-truth information). Thus, \hat{f} and \hat{X}
 68 are consistent. The minor inconsistencies in Figure 100 arise from biases existing after the
 69 convergence of the two losses.

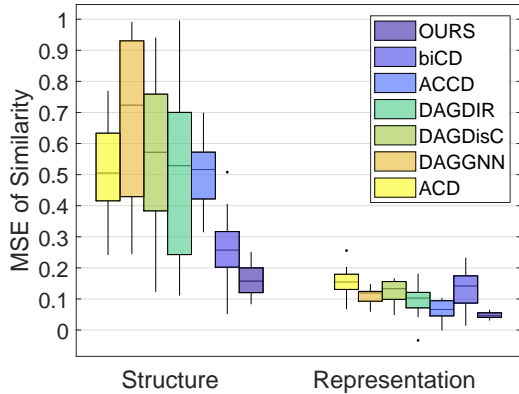


Figure 101: The boxplot showing the MSE of predicted similarity matrices to ground truth. The left clustering is the results from structure to the ground truth, and right clustering is results from representation. The similarity matrices are computed via Equation 17 and 18, respectively.

Table 100: Results of Scalability on *Causalogue* Dataset. “ $auROC_S$ ”, “ $auROC_R$ ”, and “ $auROC_C$ ” represent the AUROC of Causal Structure, Causal Representation and Causal Consistency, respectively.

Methods	20% Trainset			40% Trainset			60% Trainset			80% Trainset			100% Trainset		
	$auROC_S$	$auROC_R$	$auROC_C$	$auROC_S$	$auROC_R$	$auROC_C$	$auROC_S$	$auROC_R$	$auROC_C$	$auROC_S$	$auROC_R$	$auROC_C$	$auROC_S$	$auROC_R$	$auROC_C$
ACD	0.15	0.22	0.12	0.49	0.45	0.19	0.67	0.68	0.31	0.79	0.85	0.44	0.84	0.85	0.51
DAG-GNN	0.08	0.26	0.16	0.36	0.44	0.25	0.49	0.65	0.33	0.54	0.88	0.45	0.56	0.90	0.50
DAG-DisC	0.07	0.24	0.08	0.38	0.39	0.24	0.57	0.63	0.31	0.64	0.81	0.46	0.68	0.88	0.52
DAG-DIR	0.10	0.26	0.14	0.47	0.48	0.29	0.51	0.63	0.28	0.63	0.84	0.46	0.67	0.89	0.51
ACCD	0.13	0.19	0.12	0.45	0.46	0.27	0.61	0.65	0.39	0.74	0.87	0.46	0.79	0.93	0.60
biCD	0.16	0.25	0.18	0.53	0.49	0.26	0.74	0.69	0.37	0.84	0.82	0.57	0.91	0.86	0.64
Ours _{SSM}	0.21	0.44	0.18	0.52	0.61	0.35	0.75	0.79	0.69	0.88	0.90	0.89	0.94	0.94	0.95

70 • **Indefinite Data ($M > 1$ and $D > 1$):** For multi-structure scenarios, $\hat{f} = h_1(X, \varphi)$, and
 71 for multi-value variables, $\hat{X} = h_2(X, \hat{f})$. And $D > 1$ makes $loss(\hat{X}, X)$ ineffective.
 72 Therefore, when only $loss(\hat{f}, f)$ exists, we can get $\hat{f} = f$ and $f \Leftrightarrow X$. However, we cannot
 73 guarantee $\hat{X} \Leftrightarrow \hat{f}$ for $X = \hat{X}$, thus severe inconsistencies exist.

74 **4 VARIANCE OF LEARNING RESULTS**

75 In Table 3, we only show the evaluation results between structure (graph) and representation. To
 76 further demonstrate the benefits of our SSL framework to the model, we additionally focus on the
 77 error of the similarity matrices of structure to the ground truth, and representation to the ground truth,
 78 respectively. Figure 101 shows the box plots of the errors in structure and representation. Ours_{SSM}
 79 evidently outperforms other methods in terms of structure, even those specifically designed to handle
 80 multi-value data (ACCD, biCD). In addition, the variance of intervention-based methods (DAG-DisC,
 81 DAG-DIR) is extremely large, which aligns with our previous conclusion that intervening by negative
 82 examples leads to the additional variance of the batch size. On the representation side, almost all
 83 methods performed well. As indicated in our definition 100, nearly all methods can satisfy Correlation
 84 Condition, hence the error can be reduced to a certain extent. Nevertheless, the remaining stubborn
 85 error is due to “pseudo-correlation” caused by the inability to fully satisfy Causation Condition due
 86 to causal inconsistency. The significantly smaller variance of our method demonstrates that CCC can
 87 further help representation more completely determine Causation Condition.

88 **5 SCALABILITY**

89 We evaluate scalability by scaling the training set. Table 100 shows that our method performs best
 90 under any scale of datasets, especially in terms of structure. This is because, when the sample size is
 91 insufficient, intervention methods can extract more causal information contained in the samples. At
 92 the same time, the various real datasets in Appendix G also indirectly reflect our method’s adaptability
 93 to datasets of different scales.

94 REFERENCES

- 95 Hang Chen, Jing Luo, Xinyu Yang, and Wenjing Zhu. Affective reasoning at utterance level in
96 conversations: A causal discovery approach. *EMNLP-main*, 2023a.
- 97 Hang Chen, Xinyu Yang, and Qing Yang. Learning to recover causal relationship from indefinite
98 data in the presence of latent confounders. *arXiv preprint arXiv:2305.02640*, 2023b.
- 99 Markus Kalisch and Peter Bühlman. Estimating high-dimensional directed acyclic graphs with the
100 pc-algorithm. *Journal of Machine Learning Research*, 8(3), 2007.
- 101 Sindy Löwe, David Madras, Richard Zemel, and Max Welling. Amortized causal discovery: Learning
102 to infer causal graphs from time-series data. In *Conference on Causal Learning and Reasoning*,
103 pp. 509–525. PMLR, 2022.
- 104 Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks.
105 In *International Conference on Machine Learning*, pp. 7154–7163. PMLR, 2019.