
Multi-View Causal Discovery without Non-Gaussianity: Identifiability and Algorithms

Anonymous Authors¹

Abstract

Causal discovery is a difficult problem that typically relies on strong assumptions on the data-generating model, such as non-Gaussianity. In practice, many modern applications provide multiple related views of the same system, which has rarely been considered for causal discovery. Here, we leverage this multi-view structure to achieve causal discovery with weak assumptions. We propose a multi-view linear Structural Equation Model (SEM) that extends the well-known framework of non-Gaussian disturbances by alternatively leveraging correlation over views. We prove the identifiability of the model for acyclic SEMs. Subsequently, we propose several multi-view causal discovery algorithms, inspired by single-view algorithms (DirectLiNGAM, PairwiseLiNGAM, and ICA-LiNGAM). The new methods are validated through simulations and applications on neuroimaging data, where they enable the estimation of causal graphs between brain regions.

1. Introduction

Causal discovery is a fundamental problem in scientific data analysis as well as several technological applications (Peters et al., 2017). The basic problem is that we are given a number of observed variables, and we need to infer which variables cause which, and with what “connection strengths”. In the typical case, we only observe the data passively without any possibility of performing interventions, and this purely observational quality of the data makes the problem difficult.

Often, the problem is formalized as a structural equation model (SEM) (Bollen, 1989), also called a functional

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

causal model (Pearl, 2009). Then the question is how to estimate the SEM parameters using statistical theory. In fact, before even considering estimation methods, we need to know if it is at all possible to solve the problem. This is the question of *identifiability* of the model: can the parameters that describe the causal relations and directions be (uniquely) estimated?

A well-known fact is that if all variables are Gaussian and the model is linear, identifiability of the SEM is very problematic. Some of the earliest work showed that the directions can sometimes be recovered by an analysis of the *conditional independencies* between variables (Spirtes et al., 2001); however, this is only possible in some special cases and notably not possible when observing only two variables. Recent literature has, therefore, focused on different departures from the linear-Gaussian framework to achieve identifiability. A major advance was to consider a linear model together with *non-Gaussianity* (Shimizu et al., 2006), which leads to full identifiability of the model under weak conditions such as acyclicity; however, such non-Gaussianity may not be strong enough in many data sets. A related framework was proposed by assuming that the cause undergoes a *nonlinear transform*, while the disturbance is still additive (Hoyer et al., 2008); however the nonlinearity may not be strong enough in many applications, and such nonlinearity slightly contradicts the linear additivity of the disturbance. Further related frameworks were proposed by (Peters and Bühlmann, 2014; Zhang and Hyvärinen, 2009; Monti et al., 2019).

An alternative recent framework is given by the *multi-view* setting, where data is collected from different sources so that the data points are paired. The correlations between the views can then be leveraged for estimation and identifiability. However, current approaches using multi-view structure are scarce and limited; either they ignore dependencies across views (Shimizu, 2012) or make strong assumptions about these dependencies (Chen et al., 2024). We further note that methods in the *multi-domain* setting are fundamentally different from the multi-view case since multi-domain data is not paired: the distinction is detailed in Section 7 below.

Here, we consider the linear multi-view SEM, providing

a general theory and algorithms that leverage the multiple views. Our main contributions are:

1. We generalize three of the most well-known *algorithms* for single-view causal discovery to the multi-view setting: Pairwise LiNGAM, DirectLiNGAM and ICA-LiNGAM. The generalizations of PairwiseLiNGAM and DirectLiNGAM result in completely new and fast causal discovery using second-order statistics only, while our adaptation of ICA-LiNGAM improves on Chen et al. (2024).
2. We show that our algorithms are *consistent* under weak conditions: we can replace the non-Gaussianity assumption on the disturbances by conditions on correlations between the views and “diversity” of such Second-Order Statistics (SOS). Still, if the disturbances are non-Gaussian, we generalize Shimizu (2012); Chen et al. (2024).
3. The proofs above directly lead to rigorous *identifiability* results on the model, including cases where all the disturbances are Gaussian, or, more generally, identifiability results based on SOS only. Importantly, only weak conditions on SOS are necessary for identifiability, greatly generalizing previous work.
4. We make our algorithms usable by practitioners by providing them as *open-source*, and demonstrate their usefulness by benchmarking them on real neuroimaging data.

2. Background

Causal ordering In the following, we consider a random vector $\mathbf{x} \in \mathbb{R}^p$ whose entries are nodes of a directed acyclic graph (DAG). The DAG is represented by an adjacency matrix $\mathbf{B} \in \mathbb{R}^{p \times p}$, where non-zero values indicate the edges of the graph. The entries of \mathbf{x} can be reordered such that each “causes” the next. This is formalized by the adjacency matrix admitting a special decomposition $\mathbf{B} = \mathbf{P}^\top \mathbf{T} \mathbf{P}$, where \mathbf{T} is a strictly lower triangular matrix and \mathbf{P} is a permutation matrix referred to as the *causal ordering*. In general, the causal ordering is not unique.

LiNGAM: A single-view model for causal discovery We consider a structural equation model (SEM) called LiNGAM, that is linear and where the data follows

$$\mathbf{x} = \mathbf{B}\mathbf{x} + \mathbf{e} \quad (1)$$

and the entries e_1, \dots, e_p of the random vector $\mathbf{e} \in \mathbb{R}^p$ are independent noise terms, or *disturbances*. Causal discovery consists in inferring the parameters \mathbf{B} of the model, from observations of \mathbf{x} . Yet, the identifiability of the model

and therefore the uniqueness of \mathbf{B} is not straightforward. In fact, it is well-known that with Gaussian disturbances, the model is unidentifiable in general (Richardson and Spirtes, 2002; Genin, 2021). A major advance by Shimizu et al. (2006) was assuming *non-Gaussian* disturbances, leading to their Linear Non-Gaussian Acyclic Model (LiNGAM).

Identifiability of LiNGAM Shimizu et al. (2006) showed that the LiNGAM model is identifiable in terms of the matrix \mathbf{B} . A rigorous re-statement of this result is given in Appendix A.2. A question that has received less attention is under what conditions the model is identifiable in terms of the causal ordering \mathbf{P} . In fact, it is not in general: there may exist many permutation matrices \mathbf{P} and strictly lower triangular matrices \mathbf{T} such that the generated data has the same distribution and the generating permutation cannot be identified. For instance, in the degenerate case $\mathbf{B} = \mathbf{T} = \mathbf{0}$, any permutation matrix \mathbf{P} gives the same data distribution. As is well-known, a DAG in general defines only a *partial* ordering in the sense that for some pairs of variables, we cannot necessarily say which is “earlier” and which is “later”. Thus, if we want to make the causal ordering unique — *i.e.* to obtain a *total* ordering — we need further assumptions, as considered in later sections.

Estimation of LiNGAM Since the LiNGAM model was proposed, two important categories of algorithms for estimating this model have been developed. First, algorithms based on *recursively* computing the *residuals* of pairwise regressions include DirectLiNGAM (Shimizu et al., 2011) and PairwiseLiNGAM (Hyvärinen and Smith, 2013). These compare pairs of variables from \mathbf{x} and deduce the causal direction between them based on regressing the variables in the two directions. By aggregating this information, they finally determine a causal ordering \mathbf{P} . Once an ordering is found, the matrix \mathbf{T} can easily be obtained with conventional methods such as least-squares (LS) regression, which then gives $\mathbf{B} = \mathbf{P}^\top \mathbf{T} \mathbf{P}$.

The second category is the ICA-based algorithm that was proposed in the original LiNGAM paper (Shimizu et al., 2006). It is based on the observation that the LiNGAM model can be rewritten as a latent variable model, in particular an ICA model (Hyvärinen et al., 2001) as

$$\mathbf{x} = \mathbf{A}\mathbf{e} . \quad (2)$$

Again, the entries in \mathbf{e} are independent and non-Gaussian, and the “mixing” matrix \mathbf{A} expresses how the data is generated from the latent variable \mathbf{e} . Many methods developed for ICA can be used to estimate the matrix \mathbf{A} . However, it is important to consider the special structure of the mixing, since $\mathbf{A} = (\mathbf{I} - \mathbf{B})^{-1}$, where \mathbf{B} is a DAG (Shimizu et al., 2006). This method first recovers the adjacency matrix \mathbf{B} , and then estimates a causal ordering \mathbf{P} from \mathbf{B} .

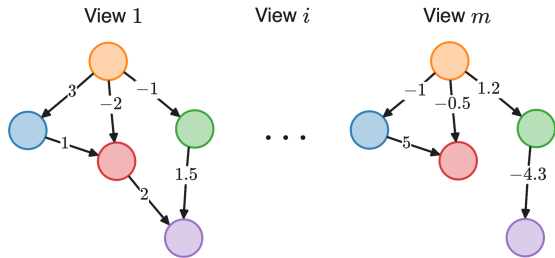


Figure 1. Visual representation of our multi-view LiMVAM model in Eq. 3. Each view of the data is described by a causal graph. The causal ordering given by the DAG is shared over views $1 \dots m$, but causal weights can differ.

3. A Multi-View Model for Causal Discovery

LiMVAM: A multi-view model for causal discovery

We now extend the single-view model in Eq. 1 to the multi-view setting, leading to a Linear Multi-View Acyclic Model (LiMVAM). Suppose we observe m different views of a p -dimensional vector. We write this as

$$\mathbf{x}^i = \mathbf{B}^i \mathbf{x}^i + \mathbf{e}^i \quad (3)$$

for views $i \in \llbracket 1, m \rrbracket$. The \mathbf{B}^i are view-specific adjacency matrices of DAGs, and \mathbf{e}^i are view-specific disturbance vectors (with entries of non-zero variance).

Similarly to Eq. 2, the LiMVAM model can equivalently be written as a latent-variable model

$$\mathbf{x}^i = \mathbf{A}^i \mathbf{e}^i \quad (4)$$

where the mixing matrix has a special structure, $\mathbf{A}^i = (\mathbf{I} - \mathbf{B}^i)^{-1}$, which we will leverage later in Section 5.

Model assumptions To exploit the multi-view structure, we follow Shimizu (2012); Chen et al. (2024) and assume that all adjacency matrices \mathbf{B}^i share the same causal ordering. Formally, they admit a decomposition $\mathbf{B}^i = \mathbf{P}^\top \mathbf{T}^i \mathbf{P}$ where \mathbf{T}^i is strictly lower triangular and \mathbf{P} is a permutation matrix that is shared among the views i . Moreover, as is commonly assumed in causal discovery, we require the components of each \mathbf{e}^i to be mutually independent, for any given i . (In contrast, correlations between views i and i' will be seen to be essential for identifiability.) Contrary to most existing approaches, we do not impose any restriction on the distribution of the disturbances, such as non-Gaussianity, which greatly broadens the applicability of the model. Figure 1 illustrates the proposed model.

We emphasize that the multi-view framework considered here is fundamentally different from the multi-domain or multi-environment settings, as it assumes that observations across different views are paired; a detailed discussion is deferred to Section 7.

The next two sections develop efficient methods for estimating the adjacency matrices \mathbf{B}^i , together with identifica-

bility theory. We begin in Section 4 with methods that estimate the causal ordering using recursive residuals and using only second-order statistics: these lead to our Pairwise-LiMVAM and DirectLiMVAM algorithms, which extend the single-view algorithms PairwiseLiNGAM (Hyvärinen and Smith, 2013) and DirectLiNGAM (Shimizu et al., 2011). In Section 5, we propose a modified version of the ICA-based multi-view model from Chen et al. (2024).

4. Estimation by SOS and recursive residuals

Here, we show how to estimate the model parameters using Second-Order Statistics (SOS) only, and how the SOS alone lead to identifiability. To do so, we will rely on a principle that we call *recursive residuals*. While it has been used by Shimizu et al. (2011); Hyvärinen and Smith (2013); Shimizu (2012), we provide here a unified treatment.

4.1. Causal ordering using recursive residuals

Methods based on what we call recursive residuals are built on the following principle: 1) we can find a root variable (*i.e.* one with no parents) by analyzing the causal directions between all the pairs of variables, and 2) when all other variables are regressed on a root variable, the resulting residuals preserve the causal ordering of the original variables (Shimizu et al., 2011, Lemma 2 and Corollary 1). In the multi-view case, the same principles apply since all views are assumed to share a common ordering. A formal proof is provided in Appendix D.3.1 and D.3.2.

Algorithms based on recursive residuals therefore follow the same recursive scheme to recover the causal ordering: (i) perform pairwise regressions to obtain residuals, (ii) identify a root variable according to a specific criterion of causal direction computed on the residuals, and (iii) remove the selected root and repeat the procedure on the residuals until all variables are ordered. In practice, we will use Ordinary LS to compute the residuals in step (iii). The main difference between algorithms lies in step (ii), *i.e.* in how the root variable is selected.

Since a root variable has no parents, we can find one by testing the causal direction for every pair of variables and taking as a root a variable that is never inferred to be an effect. (A practical way to select the root variable from such tests with finite samples is described in Appendix E.4.) Thus, root selection reduces to designing a criterion of causal direction capable of reliably distinguishing the causal direction between two variables. In the following, we introduce two new such criteria in the multi-view setting, after first discussing a condition to ensure that these criteria are consistent.

4.2. Correlation assumption (bivariate case)

Now, we consider the case of two variables, which is the foundation of the recursive residuals procedure. For simplicity of notation, we can consider these two random variables aggregated across views: $\mathbf{x}_1 = (x_1^1, \dots, x_1^m)^\top$ and $\mathbf{x}_2 = (x_2^1, \dots, x_2^m)^\top$, with either \mathbf{x}_1 causing \mathbf{x}_2 , or the converse. Covariance matrices are denoted by $\Sigma_{\mathbf{x}_1} = \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top]$, and likewise for \mathbf{x}_2 , \mathbf{e}_1 and \mathbf{e}_2 .

We introduce a fundamental assumption that makes the causal direction recoverable. Essentially, the views need to be correlated and correlations need to be diverse.

Assumption 1 (Correlation and diversity across views, bivariate case) *Consider without loss of generality that $\mathbf{x}_1 \rightarrow \mathbf{x}_2$ (in the sense that $\mathbf{B}_{21}^i \neq 0$ for at least one i). There exist two distinct views $i \neq i'$ such that a) \mathbf{x}_1 and \mathbf{x}_2 are correlated in at least one of those views, $\text{corr}(x_1^i, x_2^i) \neq 0$ or $\text{corr}(x_1^{i'}, x_2^{i'}) \neq 0$, and b) their correlations fulfill the “diversity” condition $|\text{corr}(x_1^i, x_1^{i'})| \neq |\text{corr}(e_2^i, e_2^{i'})|$.*

Discussion on the intuitive meaning of such an assumption is deferred to Section 4.6. More general conditions are discussed in Appendix B.6.

4.3. Likelihood-based criterion: PairwiseLiMVAM

We first derive a new criterion that extends the PairwiseLiNGAM criterion of Hyvärinen and Smith (2013) to the multi-view setting. Assuming Gaussian disturbances, we compute the expected log-likelihoods of two models: $\mathbf{x}_1 \rightarrow \mathbf{x}_2$ and $\mathbf{x}_2 \rightarrow \mathbf{x}_1$, where an arrow denotes a directed edge between two adjacent variables in the causal graph. This leads to the following likelihood ratio (LR) criterion:

Proposition 1 (LR criterion and its consistency) *The log-likelihood ratio of the two models is given by:*

$$\text{LR}_{12} = -\log \det \Sigma_{\mathbf{x}_1} - \log \det \Sigma_{\mathbf{e}_2} + \log \det \Sigma_{\mathbf{x}_2} + \log \det \Sigma_{\mathbf{e}_1} . \quad (5)$$

Under Assumption 1, this criterion is positive when $\mathbf{x}_1 \rightarrow \mathbf{x}_2$, and negative in the opposite direction.

We prove this in Appendix B. Note that this consistency result applies *no matter whether the disturbances are Gaussian or non-Gaussian*. Importantly, this makes the causal direction identifiable under very general circumstances. In practice, we replace the covariance matrices $\Sigma_{\mathbf{x}_1}$, $\Sigma_{\mathbf{e}_1}$, $\Sigma_{\mathbf{x}_2}$, and $\Sigma_{\mathbf{e}_2}$ with consistent estimators (see Appendix E.1).

4.4. Cross-covariance-based criterion: DirectLiMVAM

Our next criterion is a new multi-view extension of the DirectLiNGAM criterion (Shimizu et al., 2011). In a single-view setting, Shimizu et al. (2011) relied on the fact that x_1

and e_2 are *independent* if the causal direction is $x_1 \rightarrow x_2$, and measured the independence using a kernel-based estimator of mutual information called KGV (Bach and Jordan, 2002). A multi-view extension was actually proposed by Shimizu (2012), where the correct direction was chosen by comparing the *sums over views* of these scores for the two directions. However, such summation over views does not fully take advantage of the multi-view structure. Thus, the disturbances still needed to be non-Gaussian.

Here, we efficiently exploit the multi-view framework while only looking at SOS. We choose to evaluate the *cross-covariance* between the *entire* view-aggregated vectors \mathbf{x}_1 and \mathbf{e}_2 for the direction $\mathbf{x}_1 \rightarrow \mathbf{x}_2$, and between \mathbf{x}_2 and \mathbf{e}_1 for the direction $\mathbf{x}_2 \rightarrow \mathbf{x}_1$. The correct direction is chosen by comparing the Frobenius norms of the cross-covariances between regressors and residuals. Thus, we obtain the following “FC” criterion that is consistent:

Proposition 2 (FC criterion and its consistency) *Consider the following criterion:*

$$\text{FC}_{12} = \left\| \mathbb{E}[\mathbf{x}_2 \mathbf{e}_1^\top] \right\|_F - \left\| \mathbb{E}[\mathbf{x}_1 \mathbf{e}_2^\top] \right\|_F . \quad (6)$$

Under Assumption 1, this criterion is positive when $\mathbf{x}_1 \rightarrow \mathbf{x}_2$, and negative in the opposite direction.

We prove this in Appendix C.

4.5. Estimating the causal coefficients

Applying the recursive scheme explained in Section 4.1, and using the criteria from Eq. 5 or Eq. 6, we obtain a causal ordering encoded by \mathbf{P} , in the general multivariate case. Importantly, since the ordering is now known, the causal matrices \mathbf{B}^i can be uniquely recovered with little effort. In fact, the model becomes

$$\mathbf{x}^{i'} = \mathbf{T}^i \mathbf{x}^{i'} + \mathbf{e}^{i'} , \quad i \in \llbracket 1, m \rrbracket \quad (7)$$

where $\mathbf{x}^{i'} = \mathbf{P} \mathbf{x}^i$ are the reordered observations and $\mathbf{e}^{i'} = \mathbf{P} \mathbf{e}^i$ are the reordered disturbances. By construction, each variable $x_j^{i'}$ only depends on its predecessors $x_1^{i'}, \dots, x_{j-1}^{i'}$. For each $j \in \llbracket 2, p \rrbracket$, we estimate the j -th row of all matrices \mathbf{T}^i jointly using one-step Feasible Generalized Least Squares (FGLS) (Zellner, 1962), which is more efficient than Ordinary LS. A detailed description of one-step FGLS is provided in Appendix E.5. Finally, the adjacency matrices are recovered as $\mathbf{B}^i = \mathbf{P}^\top \mathbf{T}^i \mathbf{P}$.

The overall procedure is summarized in Algorithm 3 in the Appendix; it notably requires no hyperparameter tuning. Its worst-case *computational complexity* is in $O(m^3 \cdot p^3 \cdot n)$, where n is the number of observations, which we show in Appendix E.7. In practice, these algorithms are highly parallelizable as they rely on basic algebraic operations (*e.g.* multiplying pairs of matrices), and fast as they do not require iterative optimization of an objective function.

4.6. General identifiability using second-order statistics

To theoretically justify when our algorithm can find *unique* parameters \mathbf{B}^i and \mathbf{P} , we next present a general identifiability theory. It uses the following generalization of Assumption 1 from two variables to many variables.

Assumption 2 (Correlation and diversity across views) *For any two variables x_j and $x_{j'}$ such that $x_j \rightarrow x_{j'}$ (in the sense that $\mathbf{B}_{j',j}^i \neq 0$ for at least one i), there exist two distinct views $i \neq i'$ such that a) those variables are correlated in at least one of those views, $\text{corr}(x_j^i, x_{j'}^i) \neq 0$ or $\text{corr}(x_j^{i'}, x_{j'}^{i'}) \neq 0$, and b) their correlations fulfill the “diversity” condition:*

$$|\text{corr}(x_j^i, x_{j'}^{i'})| \neq |\text{corr}(e_{j'}^i, e_{j'}^{i'})|. \quad (8)$$

Even milder sufficient conditions are discussed in Appendix D.1 but they are harder to interpret. To better understand this assumption, note first that it requires some *correlation across views*: it rules out the trivial case where all cross-view correlations are zero. Moreover, it demands some *diversity* in these correlations. Consider the extreme situation where, for standardized variables, $x_j^i = \alpha x_{j'}^{i'}$ and $x_{j'}^i = \alpha x_j^{i'}$ (and analogously for the residuals), for some scalar α . In this case, all correlations equal 1, so Eq. 8 is violated. Appendix D.2 provides more realistic examples of violations. We next introduce an additional assumption that ensures identifiability of the total causal ordering.

Assumption 3 (Dense connectivity of the graph when pooled across views) *Let \mathcal{G} denote the union graph of the m views, obtained by taking the union of their edge sets. Assume that the ordering induced by \mathcal{G} is total, i.e. there exists a directed path between any two variables.*

With these assumptions, we prove in Appendix D.3 the following identifiability result.

Theorem 3 (Identifiability of the ordering and adjacency matrices \mathbf{B}^i) *Assume the data follows the model in Eq. 3, and that Assumption 2 holds. Then, the causal (partial) ordering and causal coefficients \mathbf{B}^i are identifiable. If we further make Assumption 3, then \mathbf{P} is identifiable as well.*

The identifiability of the *partial* ordering does not necessarily mean that the permutation matrix \mathbf{P} can be recovered uniquely. If the ordering is only partial (for example, in a three-nodes graph: $x_1 \rightarrow x_2$ and $x_1 \rightarrow x_3$, but there is no relation between x_2 and x_3 , so both orderings $1 \rightarrow 2 \rightarrow 3$ and $1 \rightarrow 3 \rightarrow 2$ are valid), then \mathbf{P} is only identifiable up to a class of permutations that induce the same partial ordering. The above theorem shows a case where \mathbf{P} can be (uniquely) identified.

5. Estimation by multi-view ICA assuming shared disturbances

Next, we provide an alternative approach to the LiM-VAM model definition and estimation. We utilize possible non-Gaussianity of the disturbances to obtain identifiability conditions that are different from the previous section. However, we still do not impose non-Gaussianity as a necessary condition.

Model definition The approach here is to propose a particular model for the disturbances by supposing that they can be additively composed into two terms: one term is view-specific, while the other is shared across views and therefore makes the views dependent. Specifically:

Assumption 4 (Shared disturbances) *The disturbances can be decomposed as*

$$e^i = \mathbf{D}^i \mathbf{s} + \mathbf{n}^i \quad (9)$$

where the \mathbf{D}^i are diagonal matrices with positive entries on the diagonal, the \mathbf{s} has mutually independent entries and second-order moment $\mathbb{E}[\mathbf{s}\mathbf{s}^\top] = \mathbf{I}_p$, and the view-specific noises are Gaussian $\mathbf{n}^i \sim \mathcal{N}(\mathbf{0}, \Sigma^i)$ and depend on the view i via the second-order moments Σ^i that are diagonal matrices. Lastly, the vectors \mathbf{s} and $\mathbf{n}^i, i = 1, \dots, m$, are assumed to be mutually independent.

This assumption and Eq. 9 are related to Chen et al. (2024) who were originally inspired by the multi-view Shared ICA of Richard et al. (2021). Specifically, Chen et al. (2024) considered the special case where the matrices $\mathbf{D}^i = \mathbf{I}_p$, so that there are the same disturbances \mathbf{s} over views. We argue here that their model was too restrictive and unrealistic as it imposes an arbitrary scaling for the data variables. Consider the root variable of the DAG, denoting it by x_i . Its disturbance is $s_i + n_i$ which has variance of at least one since the variance of s_i is defined as one. Thus, the root variable has variance of at least one which is absurd: the model should be invariant to the scaling of the variables. Likewise, suppose we create a new data set by dividing just one view, say x^1 by, say 10, giving new data x'^1 . This new data does not follow the model by Chen et al. (2024) anymore since \mathbf{s} should be rescaled to $\mathbf{s}/10$ for just this view; now the $\mathbf{s}/10$ here is different from the \mathbf{s} in other views. Such problems motivated us to develop a model where the scaling terms \mathbf{D}^i are new parameters, and thus to the definition in Eq. 9.

Identifiability Now we proceed to show the identifiability of the model based on Assumption 4. The starting point is that our model can be rewritten as a multi-view ICA model as in Eq. 4, just like in the case of LiNGAM. Thus, the methods developed for multi-view ICA can be used to

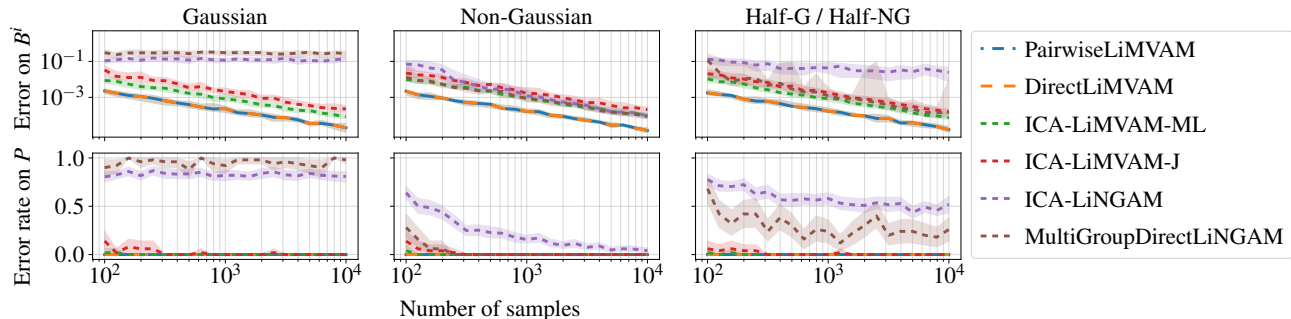


Figure 2. Separation performance of three recursive residuals algorithms and three ICA-based algorithms. We assume shared disturbances as in Eq. 9. We used 50 different seeds. Lower is better.

analyze its identifiability. In the case of Gaussian disturbances, we use the following assumption.

Assumption 5 (Diversity of the views with Gaussian disturbances) *If j and j' , $j \neq j'$, are the indices of two Gaussian common disturbances in s , then the number of views is at least three, and there exists a view i such that*

$$\frac{\Sigma_{jj}^i}{(D_{jj}^i)^2} \neq \frac{\Sigma_{j'j'}^i}{(D_{j'j'}^i)^2}.$$

This assumption, inspired by Richard et al. (2021), is trivially fulfilled if all disturbances are non-Gaussian. For Gaussian disturbances, it essentially states that the views should be diverse in terms of disturbances’ SOS. This follows from the theory of Shared ICA which does not require the common disturbances s to be non-Gaussian if there is enough diversity (Richard et al., 2021; Anderson et al., 2013).

We can now state our main identifiability result for this variant of the model.

Theorem 4 (Identifiability with shared disturbances) *Assume the LiMVAM model in Eq. 3 together with Assumption 4. Under Assumption 5, the model is identifiable in the sense that the parameters B^i , D^i , and Σ^i are identifiable. If we further make Assumption 3, P is identifiable as well.*

The proof of the theorem is given in Appendix F. We thus see that for non-Gaussian disturbances, identifiability is achieved with weak conditions; in fact, no special conditions on the covariances are needed. Still, the theory of Shared ICA leads to identifiability even for Gaussian disturbances, but with stronger conditions than obtained in Section 4. Moreover, the Shared ICA framework naturally provides an algorithm.

Estimation The model can be estimated by a combination of the ICA-based estimation procedure by Shimizu et al. (2006), combined with the Shared ICA methods by Richard et al. (2021). We call it ICA-LiMVAM and explain it in detail in Algorithm 4. Surprisingly, the resulting algorithm gives estimates of B_i that are identical to Chen

et al. (2024), even though our model is more general. In contrast, our algorithm estimates different noise variances and scalings D_i . The worst-case *computational complexity* is in $O(tnmp^3 + mp^5)$, for t iterations of the Shared ICA (ML) algorithm, n samples, m views and p components. We detail this in Appendix G.2.

6. Experiments

We next apply our algorithms to synthetic data and real-world neuroimaging data; code will be released publicly.

6.1. Simulations

We benchmark a total of six multi-view causal discovery algorithms on synthetic data. Three of these algorithms are based on recursive residuals: our PairwiseLiMVAM and DirectLiMVAM, as well as MultiGroupDirectLiNGAM from Shimizu (2012). The other three algorithms are ICA-based: ICA-LiMVAM which is essentially the same as Chen et al. (2024), as well as ICA-LiNGAM (Shimizu et al., 2006) which is naively applied to each view, separately. Note that ICA-LiMVAM comes in two variants, each using a different ICA algorithms: either Shared ICA “ML”, a likelihood-based method that can handle both Gaussian and non-Gaussian disturbances, or Shared ICA “J”, which jointly diagonalizes covariance matrices without using non-Gaussianity. We do not include comparisons with multi-domain causal discovery methods, as these are not directly comparable to multi-view approaches in settings with cross-view correlations (see Section 7 for a detailed discussion). Nevertheless, for completeness, we provide an additional comparison between DirectLiMVAM and a multi-domain method (Perry et al., 2022) in Appendix H.1.3, where we see that our DirectLiMVAM achieves better performance.

The first synthetic experiment is inspired by Richard et al. (2021). We monitor the performance of causal discovery algorithms across varying sample sizes and disturbance distributions. The data are generated according to the LiM-

VAM model with shared disturbances from Eq. 3 and Eq. 9. The performance of the causal discovery algorithms is measured by the ℓ_2 distance between true and estimated adjacency matrices B^i , and the percentage of incorrectly estimated causal orders P over 50 random runs.

Results of the first synthetic experiment are plotted in Figure 2. Our algorithms PairwiseLiMVAM and DirectLiMVAM consistently outperform all others by a significant margin. ICA-LiMVAM yields reliable estimates but with larger errors. In contrast, MultiGroupDirectLiNGAM and ICA-LiNGAM, which are designed for non-Gaussian disturbances, fail entirely in the Gaussian setting, as expected. Note that the high error rate of ICA-LiNGAM in recovering the causal ordering P is likely explained by its inability to exploit that the ordering is shared across views.

For the second synthetic experiment, we report the trade-off between computational complexity and estimation error of the adjacency matrices. We use a different data generation model, described by Eq. 3 where the disturbances are Gaussian with equal variances, and without Eq. 9. This breaks the assumptions for all methods but our PairwiseLiMVAM and DirectLiMVAM. As expected, we see in Figure 3 that our algorithms PairwiseLiMVAM and DirectLiMVAM provide a major reduction in estimation errors, at little or no computational cost, while the four other algorithms struggle.

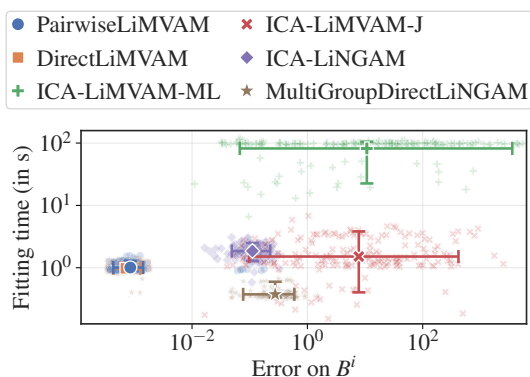


Figure 3. Estimation error of the adjacency matrices B^i vs. total fitting time for each method, across 200 random runs with $m = 6$ views, $p = 5$ variables, and $n = 1000$ samples. Disturbances are drawn from a standard Gaussian. Note that the point clouds for PairwiseLiMVAM and DirectLiMVAM largely overlap.

In Appendix H.1, we include further experiments. They test Assumption 5 and the effect of the adjacency matrices’ sparsity, as well as scalability with more views and variables. In particular, we find that DirectLiMVAM outperforms PairwiseLiMVAM in some cases, which is less obvious from Figures 2 and 3 alone.

6.2. Real data experiments with MEG

Next, we performed experiments on magnetoencephalography (MEG) data measuring human brain activity. MEG data has a high temporal resolution, so typically methods related to Granger causality are used. However, we are here interested in analyzing the *energies of oscillatory* signals. The energies change very slowly (on the time-scale of seconds), while the underlying brain activity being measured is likely to change much more quickly. Thus, it is appropriate to use a model of instantaneous causality like ours.

We used the Cam-CAN dataset, which is the largest publicly available MEG dataset (Shafto et al., 2014; Taylor et al., 2017), and considered the “sensorimotor task” during which each participant had to respond with a right index finger button press to auditory/visual stimuli. Considering each of the 98 participants to be a view, we applied our PairwiseLiMVAM. The experiment ran in 2 hours and 20 minutes using 5 CPUs. One may reasonably hypothesize that the participants’ brains share structural patterns, so to facilitate population-level interpretation, we computed the element-wise median of the B^i matrices.

Fig. 4a displays the ten strongest (in absolute value) median causal effects. Notably, many causal connections are between homologous regions across hemispheres — particularly within the primary motor and somatosensory cortices — consistent with prior findings that symmetric brain regions exhibit strong inter-hemispheric correlations (Li et al., 1996). Additionally, numerous arrows involve the postcentral regions, which have been implicated in fine motor control and hand-related tasks (Braun et al., 2002). To evaluate the robustness of the results, we repeated the experiment 50 times, each time randomly selecting 30 participants from the full cohort. We assessed consistency by measuring the Pearson correlation between the resulting median matrices across runs. As shown in Figure 4b, the median effects are highly correlated across different subsets (average correlation: 0.67), demonstrating the stability of our method. Additional experiments using ICA-LiMVAM-ML are reported in Appendix H.2.2.

6.3. Real data experiments with fMRI

We also evaluated PairwiseLiMVAM on an fMRI rhyming judgment task dataset (Ramsey et al., 2010). The data consist of recordings from nine participants, each represented by nine variables: one task regressor (Input I), obtained by convolving the boxcar design of the rhyming task with a canonical hemodynamic response function, and eight regional time series from bilateral regions of interest. The dataset was acquired on a 3T scanner with a repetition time of 2 s, yielding 160 samples per subject, and is available



Figure 4. (a) Top ten strongest median causal effects (estimated by PairwiseLiMVAM) across 98 Cam-CAN (MEG) participants. (b) Pearson correlations between median causal matrices across 50 runs in MEG. (c) Top ten strongest median causal effects from PairwiseLiMVAM across 9 fMRI participants. Arrows indicate direction and strength of effects, with width proportional to magnitude.

via the OpenfMRI project (preprocessed version from ¹).

We report the ten strongest effects from the element-wise median of the estimated adjacency matrices in Figure 4c. The recovered structure is consistent with established findings for this task. In particular, we find a strong edge from the task regressor to the left occipital cortex (Input \rightarrow LOCC), reflecting the expected visual drive of the stimuli, and prominent connections from left to right homologous regions (*e.g.* LOCC \rightarrow ROCC, LIFG \rightarrow RIFG, LIPL \rightarrow RIPL), in line with the left-hemisphere dominance and inter-hemispheric influences reported by Ramsey et al. (2010). Edges such as LOCC \rightarrow LIPL and LIPL \rightarrow LIFG, paths such as LOCC \rightarrow LACC, and connections between LACC and RACC also match relationships highlighted in their analyses.

7. Discussion and related work

Multi-view vs. multi-domain distinction We wish to emphasize that the multi-view setting considered here is very different from *multi-domain* methods. As an intuitive example, multi-view would correspond to a person seeing the same object from different angles, while multi-domain would correspond to a person seeing different objects. In multi-view data, the data points in different views are *paired* over domains, while in multi-domain methods, no such correspondence between data points exists. In statistical terms, the fundamental difference is that in the multi-view setting, the views are statistically dependent (correlated), while in the multi-domain setting, the domains are independently generated, typically with a distribution shift over domains. In fact, several methods study causal discovery in the multi-domain case (Ghassami et al., 2018; Adams et al., 2021; Mooij et al., 2020); *multi-environment* methods are a special case of multi-domain methods, where distributional shifts arise from interventions (Peters et al., 2016; Perry et al., 2022). Appendix I discusses those methods in more detail.

¹<https://github.com/cabal-cmu/Feedback-Discovery>

Related work A multi-view causal model based on LiNGAM was considered by Shimizu (2012), but their model was very different. The fundamental difference is that we assume the disturbances have some shared information across views and can be Gaussian, whereas Shimizu (2012) estimate necessarily non-Gaussian disturbances, while ignoring dependencies across views. Thus, they only showed how to use the multiple views to improve the estimation while still assuming non-Gaussianity of the disturbances.

Chen et al. (2024) considered multi-view SEM using a special case of our shared disturbances model, which is in its turn a special case of our general LiMVAM model. Algorithmically, there is little difference between their work and our ICA-LiMVAM. However, their identifiability proofs are not available nor did they consider the identifiability of the causal ordering, while we do both. Moreover, as we already argued in Section 5, their model formulation has the serious flaw that it forces an arbitrary scaling of the data variables. Crucially, we propose new algorithms, PairwiseLiMVAM and DirectLiMVAM, that use only SOS; they are shown to be empirically superior to the ICA-LiMVAM in most cases.

8. Conclusion

This work provides a full treatment of multi-view causal discovery that was under-explored in the literature. In particular, we define a multi-view model called LiMVAM. We derive a comprehensive identifiability theory of the causal ordering and weights, using second-order statistics (SOS, *i.e.* covariance structure) only. Beyond these theoretical guarantees, we also improve an existing algorithm (ICA-LiMVAM) and introduce two novel ones — PairwiseLiMVAM and DirectLiMVAM — for estimating the parameters of the multi-view model; again, SOS alone are enough. We rigorously prove the consistency of our algorithms and empirically show that they are both *fast* and result in *highly efficient estimation*. Experiments on synthetic data show that they outperform comparable algorithms in many settings, and results on brain imaging data are promising.

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- J. Adams, N. Hansen, and K. Zhang. Identification of partially observed linear causal models: Graphical conditions for the non-gaussian and heterogeneous cases. *Advances in Neural Information Processing Systems*, 34: 22822–22833, 2021.
- M. Anderson, G. Fu, R. Phlypo, and T. Adali. Independent vector analysis: Identification conditions and performance bounds. *IEEE Transactions on Signal Processing*, 62:4399–4410, 2013.
- F. R. Bach and M. I. Jordan. Kernel independent component analysis. *Journal of machine learning research*, 3 (Jul):1–48, 2002.
- K. A. Bollen. *Structural equations with latent variables*, volume 210. John Wiley & Sons, 1989.
- C. Braun, M. Haug, K. Wiech, N. Birbaumer, T. Elbert, and L. E. Roberts. Functional organization of primary somatosensory cortex depends on the focus of attention. *Neuroimage*, 17(3):1451–1458, 2002.
- W. Chen, X. Huang, Z. Li, R. Cai, Z. Huang, and Z. Hao. Individual causal structure learning from population data. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 7109–7117, 2024.
- P. Comon. Independent component analysis, a new concept? *Signal Processing*, 36(3):287–314, 1994. ISSN 0165-1684. Higher Order Statistics.
- A. M. Dale, B. Fischl, and M. I. Sereno. Cortical surface-based analysis: I. segmentation and surface reconstruction. *Neuroimage*, 9(2):179–194, 1999.
- K. Genin. Statistical undecidability in linear, non-gaussian causal models in the presence of latent confounders. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 13564–13574. Curran Associates, Inc., 2021.
- A. Ghassemi, N. Kiyavash, B. Huang, and K. Zhang. Multi-domain causal structure learning in linear systems. *Advances in neural information processing systems*, 31, 2018.
- A. Gramfort, M. Luessi, E. Larson, D. A. Engemann, D. Strohmeier, C. Brodbeck, R. Goj, M. Jas, T. Brooks, L. Parkkonen, and M. Hämäläinen. Meg and eeg data analysis with mne-python. *Frontiers in Neuroscience*, 7, 2013. ISSN 1662-453X.
- A. Gramfort, M. Luessi, E. Larson, D. A. Engemann, D. Strohmeier, C. Brodbeck, L. Parkkonen, and M. S. Hämäläinen. Mne software for processing meg and eeg data. *NeuroImage*, 86:446–460, 2014. ISSN 1053-8119.
- W. H. Greene. *Econometric analysis*. Pearson education india, 2003.
- M. S. Hämäläinen and R. J. Ilmoniemi. Interpreting magnetic fields of the brain: minimum norm estimates. *Medical & Biological Engineering & Computing*, 32(1):35–42, 1994.
- P. Hoyer, D. Janzing, J. M. Mooij, J. Peters, and B. Schölkopf. Nonlinear causal discovery with additive noise models. *Advances in neural information processing systems*, 21, 2008.
- A. Hyvärinen and S. M. Smith. Pairwise likelihood ratios for estimation of non-Gaussian structural equation models. *Journal of Machine Learning Research*, 14:111–152, 2013.
- A. Hyvärinen, J. Karhunen, and E. Oja. *Independent Component Analysis*. Wiley Interscience, 2001.
- S. Khan, J. A. Hashmi, F. Mamashli, K. Michmizos, M. G. Kitzbichler, H. Bharadwaj, Y. Bekhti, S. Ganesan, K.-L. A. Garel, S. Whitfield-Gabrieli, R. L. Gollub, J. Kong, L. M. Vaina, K. D. Rana, S. M. Stuffelbeam, M. S. Hämäläinen, and T. Kenet. Maturation trajectories of cortical resting-state networks depend on the mediating frequency band. *NeuroImage*, 174:57–68, 2018. ISSN 1053-8119.
- A. Li, F. Z. Yetkin, R. Cox, and V. M. Haughton. Ipsilateral hemisphere activation during motor and sensory tasks. *American Journal of Neuroradiology*, 17(4):651–655, 1996. ISSN 0195-6108.
- R. P. Monti, K. Zhang, and A. Hyvärinen. Causal discovery with general non-linear relationships using non-linear ICA. In *Proc. 35th Conf. on Uncertainty in Artificial Intelligence (UAI2019)*, Tel Aviv, Israel, 2019.
- J. M. Mooij, S. Magliacane, and T. Claassen. Joint causal inference from multiple contexts. *Journal of machine learning research*, 21(99):1–108, 2020.
- J. Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, USA, 2nd edition, 2009. ISBN 052189560X.

- 495 R. Perry, J. Von Kügelgen, and B. Schölkopf. Causal dis-
496 covery in heterogeneous environments under the sparse
497 mechanism shift hypothesis. *Advances in Neural Infor-*
498 *mation Processing Systems*, 35:10904–10917, 2022.
- 499 J. Peters and P. Bühlmann. Identifiability of gaussian
500 structural equation models with equal error variances.
501 *Biometrika*, 101(1):219–228, 2014.
- 502 J. Peters, P. Bühlmann, and N. Meinshausen. Causal in-
503 ference by using invariant prediction: identification and
504 confidence intervals. *Journal of the Royal Statistical So-*
505 *ciety Series B: Statistical Methodology*, 78(5):947–1012,
506 2016.
- 507 J. Peters, D. Janzing, and B. Schölkopf. *Elements of causal*
508 *inference: foundations and learning algorithms*. The
509 MIT Press, 2017.
- 510 L. Power, C. Allain, T. Moreau, A. Gramfort, and T. Bar-
511 douille. Using convolutional dictionary learning to de-
512 tect task-related neuromagnetic transients and ageing
513 trends in a large open-access dataset. *NeuroImage*, 267:
514 119809, 2023. ISSN 1053-8119.
- 515 J. D. Ramsey, S. J. Hanson, C. Hanson, Y. O. Halchenko,
516 R. A. Poldrack, and C. Glymour. Six problems for
517 causal inference from fmri. *neuroimage*, 49(2):1545–
518 1558, 2010.
- 519 H. Richard, P. Ablin, B. Thirion, A. Gramfort, and A. Hy-
520 varinen. Shared independent component analysis for
521 multi-subject neuroimaging. In *Advances in Neural In-*
522 *formation Processing Systems*, volume 34, pages 29962–
523 29971. Curran Associates, Inc., 2021.
- 524 T. Richardson and P. Spirtes. Ancestral graph Markov mod-
525 els. *The Annals of Statistics*, 30(4):962 – 1030, 2002.
- 526 M. A. Shafto, L. K. Tyler, M. Dixon, J. R. Taylor,
527 J. B. Rowe, R. Cusack, A. J. Calder, W. D. Marslen-
528 Wilson, J. S. Duncan, T. Dalgleish, R. N. A. Henson,
529 C. Brayne, and F. E. Matthews. The cambridge centre
530 for ageing and neuroscience (cam-can) study protocol: a
531 cross-sectional, lifespan, multidisciplinary examination
532 of healthy cognitive ageing. *BMC Neurology*, 14, 2014.
- 533 S. Shimizu. Joint estimation of linear non-gaussian acyclic
534 models. *Neurocomputing*, 81:104–107, 2012.
- 535 S. Shimizu. *Statistical Causal Discovery: LiNGAM Ap-*
536 *proach*. Springer, 2022.
- 537 S. Shimizu, P. O. Hoyer, A. Hyvärinen, A. Kerminen,
538 and M. Jordan. A linear non-gaussian acyclic model
539 for causal discovery. *Journal of Machine Learning Re-*
540 *search*, 7(10), 2006.
- 541 S. Shimizu, T. Inazumi, Y. Sogawa, A. Hyvarinen,
542 Y. Kawahara, T. Washio, P. O. Hoyer, K. Bollen, and
543 P. Hoyer. Directlingam: A direct method for learning a
544 linear non-gaussian structural equation model. *Journal*
545 *of Machine Learning Research-JMLR*, 12(Apr):1225–
546 1248, 2011.
- 547 P. Spirtes, C. Glymour, and R. Scheines. *Causation, pre-*
548 *dition, and search*. MIT press, second edition edition,
549 2001.
- 549 N. Sturma, C. Squires, M. Drton, and C. Uhler. Unpaired
550 multi-domain causal representation learning. In A. Oh,
551 T. Naumann, A. Globerson, K. Saenko, M. Hardt, and
552 S. Levine, editors, *Advances in Neural Information Pro-*
553 *cessing Systems*, volume 36, pages 34465–34492. Cur-
554 ran Associates, Inc., 2023.
- 555 S. Taulu and J. Simola. Spatiotemporal signal space sep-
556 aration method for rejecting nearby interference in meg
557 measurements. *Physics in Medicine & Biology*, 51(7):
558 1759, mar 2006.
- 559 J. R. Taylor, N. Williams, R. Cusack, T. Auer, M. A. Shafto,
560 M. Dixon, L. K. Tyler, Cam-CAN, and R. N. Henson.
561 The cambridge centre for ageing and neuroscience (cam-
562 can) data repository: Structural and functional mri, meg,
563 and cognitive data from a cross-sectional adult lifespan
564 sample. *NeuroImage*, 144:262–269, 2017. ISSN 1053-
565 8119. Data Sharing Part II.
- 566 A. Zellner. An efficient method of estimating seem-
567 ingly unrelated regressions and tests for aggregation
568 bias. *Journal of the American Statistical Association*,
569 57(298):348–368, 1962. doi: 10.2307/2281644.
- 570 K. Zhang and A. Hyvärinen. On the identifiability of
571 the post-nonlinear causal model. In *Proc. 25th Confer-*
572 *ence on Uncertainty in Artificial Intelligence (UAI2009)*,
573 pages 647–655, Montréal, Canada, 2009.

Appendix

The appendix is organized as follows.

In Section A, we rigorously restate known identifiability results for LiNGAM.

In Section B, we prove the consistency of the likelihood-based criterion for a pair of variables.

In Section C, we prove the consistency of the cross-covariance-based criterion for a pair of variables.

In Section D, we extend the consistency of these criteria for a pair of variables to *many* pairs, and show that the LiMVAM model is identifiable.

In Section E, we show how to implement algorithms used to identify the causal graph.

In Section F, we prove the identifiability of a special case of LiMVAM with shared disturbances.

In Section G, we present an algorithm for estimating the causal graph in this special case.

In Section H, we detail our experiments on real and synthetic data.

In Section I, we discuss the difference between multi-view and multi-domain frameworks.

A	LiNGAM	13
A.1	Technical lemmas	13
A.2	Identifiability of LiNGAM	13
A.3	Proofs of technical lemmas	14
B	Likelihood-based criterion	16
B.1	Basic setting	16
B.2	Derivation of the expected log-likelihood in one direction	16
B.3	Derivation of the likelihood-based criterion	16
B.4	Non-negativity of the likelihood-based criterion	17
B.5	Proof of a useful lemma	18
B.6	Conditions for strict positivity of the likelihood-based criterion	19
C	Cross-covariance-based criterion	22
C.1	Derivation of the cross-covariance-based criterion	22
C.2	Conditions for strict positivity of the cross-covariance criterion	22
D	Identifiability of LiMVAM	23
D.1	Multivariate conditions for strict positivity of a criterion	23
D.2	An example for interpreting the SOS-based identifiability condition in Assumption 2	23
D.3	Proof of Theorem 3	25
E	SOS-based Algorithms for LiMVAM	27
E.1	Consistent estimator of the likelihood-based criterion	27
E.2	Consistent estimator of the cross-covariance-based criterion	27
E.3	Algorithms for estimating the two criteria	28

605	E.4 How to find the root variable	28
606	E.5 One-step Feasible Generalized Least Squares (FGLS)	29
607	E.6 Full algorithm	30
608	E.7 Computational complexity	30
609		
610		
611		
612	F Identifiability of LiMVAM with shared disturbances	32
613	F.1 Proof of Theorem 4 (first claim)	32
614	F.2 Proof of Theorem 4 (second claim)	32
615	F.3 Results for view-specific causal orderings	33
616		
617		
618		
619	G ICA-based Algorithm for LiMVAM with shared disturbances	35
620	G.1 ICA-based Algorithm	35
621	G.2 Computational complexity	35
622		
623		
624	H Experiments	37
625	H.1 Synthetic experiments	37
626	H.2 Real data experiments	40
627		
628		
629	I Distinction between multi-view and multi-domain frameworks	42
630		
631		
632		
633		
634		
635		
636		
637		
638		
639		
640		
641		
642		
643		
644		
645		
646		
647		
648		
649		
650		
651		
652		
653		
654		
655		
656		
657		
658		
659		

660 A. LiNGAM

661 **Defining a domain connecting ICA and SEM parameters** In the following proofs, we define the set

$$662 \mathcal{W} = \{W \in \mathbb{R}^{p \times p} \mid \text{there exist a diagonal matrix } D \text{ with positive diagonal entries} \\ 663 \text{ and a DAG matrix } B, \text{ such that } W = D^{-1}(I - B)\} \quad (10)$$

664 as the domain of the “unmixing matrices” $W = A^{-1}$ in the ICA model from Eq. 2, that are compatible with the DAG
665 structure in the structural equation model (SEM) from Eq. 1. In general, unmixing matrices are not constrained in the ICA
666 model, except by invertibility (and possibly by normalizing their rows). However, in the context of SEM estimation, the
667 requirement that $W = D^{-1}(I - B)$ belongs to the set \mathcal{W} is a key structural constraint. Furthermore, it should be noted
668 that finding W is equivalent to finding D and B , since they are related by $B = D^{-1}(I - W)$.

669 Finally, we will use the fact that any DAG matrix B can be decomposed into $B = P^\top T P$, where P is a permutation
670 matrix and T a strictly lower triangular matrix.

671 A.1. Technical lemmas

672 The following two lemmas are about matrices in the context of DAGs. They are used in the proofs of Theorems 7, 4,
673 and 10, and are proven in Appendix A.3.

674 The basic identifiability results for the causal matrices B^i use the following:

675 **Lemma 5** *Let $W = D^{-1}(I - B)$ and $W' = D'^{-1}(I - B')$ be two matrices that belong to \mathcal{W} (defined in Eq. 10), and
676 let Q be a sign-permutation matrix. Then*

$$677 W' = Q^\top W \implies Q = I, D' = D, \text{ and } B' = B. \quad (11)$$

678 On the other hand, the identifiability results for the causal ordering P (or causal orderings P^i) use the following:

679 **Lemma 6** *Let P and P' be permutation matrices and T and T' be strictly lower triangular matrices, such that $P^\top T P =$
680 $P'^\top T' P'$. Assume that T contains only non-zero elements in its strictly lower triangular part. Then, $P = P'$ and
681 $T = T'$.*

682 A.2. Identifiability of LiNGAM

683 As a special case, we restate here the identifiability of the single-view LiNGAM model (in terms of the matrix B).

684 **Theorem 7** (Identifiability of LiNGAM) *In the statistical model defined by Eq. 1, the parameter B is identifiable, provided
685 that the entries in s are mutually independent, that at most one of them is Gaussian, and that B is a DAG.*

686 This proof is based on the one from Shimizu et al. (2006), which we attempt to make a bit more rigorous. Identifiability
687 was also shown in Shimizu (2022, Section 2.3) in the case with 2 components.

688 Let us transform the SEM

$$689 x = Bx + s \quad (12)$$

690 into an ICA model

$$691 x = As \quad (13)$$

692 where the mixing matrix A is constrained as $A^{-1} = I - B$, and where B is a DAG. The assumptions on s are identical
693 in SEM and ICA. Let us parameterize by $W = A^{-1} = I - B$, instead of B .

694 Consider two matrices $W = I - B$ and $W' = I - B'$, with B and B' being DAG matrices, that parameterize the same
695 statistical model in Eq. 12. We have $\det(W) = \det(W') = 1$, so W and W' are invertible, which makes them valid
696 unmixing matrices for the same ICA model in Eq. 13. So, from the identifiability theory of ICA (Comon, 1994), we know
697 that there exist a sign-permutation matrix Q and a scaling matrix D (i.e. a diagonal matrix with positive entries on the
698 diagonal) such that

$$699 W' = DQ^\top W \quad (14)$$

715 hence

$$716 \quad \mathbf{W}'' = \mathbf{Q}^\top \mathbf{W} \quad (15)$$

717 where we defined $\mathbf{W}'' = \mathbf{D}^{-1} \mathbf{W}'$. We thus have that $\mathbf{W}'' = \mathbf{D}^{-1}(\mathbf{I} - \mathbf{B}')$ and $\mathbf{W} = \mathbf{I} - \mathbf{B}$ both belong to the set \mathcal{W} ,
 718 so applying Lemma 5 to \mathbf{W} and \mathbf{W}'' imposes $\mathbf{Q} = \mathbf{I}$, $\mathbf{D} = \mathbf{I}$, and $\mathbf{B}' = \mathbf{B}$. This makes \mathbf{B} fully identifiable (“fully”
 719 identifiable is in contrast to conventional ICA theory where some indeterminacies always remain), which concludes the
 720 proof.
 721

722 A.3. Proofs of technical lemmas

723 A.3.1. PROOF OF LEMMA 5

724 Consider $\mathbf{W}, \mathbf{W}' \in \mathcal{W}$ and \mathbf{Q} a sign-permutation matrix. Suppose that

$$725 \quad \mathbf{W}' = \mathbf{Q}^\top \mathbf{W} . \quad (16)$$

726 **A permutation inequality** Using the definition of \mathcal{W} in Eq. 10 and the decompositions of DAG matrices, we have

$$727 \quad \mathbf{W} = \mathbf{D}^{-1} \mathbf{P}^\top (\mathbf{I} - \mathbf{T}) \mathbf{P} , \quad \mathbf{W}' = \mathbf{D}'^{-1} \mathbf{P}'^\top (\mathbf{I} - \mathbf{T}') \mathbf{P}' \quad (17)$$

728 where \mathbf{D}, \mathbf{D}' are diagonal matrices with positive entries on the diagonal, \mathbf{P}, \mathbf{P}' are permutation matrices, and \mathbf{T}, \mathbf{T}' are
 729 strictly lower triangular matrices. Denote $\mathbf{L} = \mathbf{I} - \mathbf{T}$ and $\mathbf{L}' = \mathbf{I} - \mathbf{T}'$; these are lower triangular matrices that have unit
 730 diagonals. We plug the decompositions into Eq. 16, and get

$$731 \quad \mathbf{D}'^{-1} \mathbf{P}'^\top \mathbf{L}' \mathbf{P}' = \mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{P}^\top \mathbf{L} \mathbf{P} . \quad (18)$$

732 To exploit the structure of the lower triangular matrices \mathbf{L} and \mathbf{L}' and show how they constrain \mathbf{Q} , we now switch notations
 733 from permutation — or sign-permutation — matrices $(\mathbf{P}, \mathbf{P}', \mathbf{Q})$, to their corresponding permutation functions (ϕ, ϕ', ψ) .
 734 Eq. 18 thus yields

$$735 \quad \frac{\mathbf{L}'_{\phi'(i), \phi'(j)}}{\mathbf{D}'_{ii}} = \pm \frac{\mathbf{L}_{\phi(\psi(i)), \phi(j)}}{\mathbf{D}_{\psi(i), \psi(i)}} , \quad \forall i, j \in \llbracket 1, p \rrbracket \quad (19)$$

736 where the \pm symbol comes from \mathbf{Q} . In particular, \mathbf{L}' has unit diagonal, so

$$737 \quad \mathbf{L}_{\phi(\psi(i)), \phi(i)} = \pm \frac{\mathbf{D}_{\psi(i), \psi(i)}}{\mathbf{D}'_{ii}} \neq 0 , \quad \forall i \in \llbracket 1, p \rrbracket \quad (20)$$

738 and \mathbf{L} is lower triangular, so that its non-zero entries must satisfy

$$739 \quad \phi(i) \leq \phi(\psi(i)) , \quad \forall i \in \llbracket 1, p \rrbracket . \quad (21)$$

740 More generally, we can replace i with $\psi(i)$, and so on, and obtain

$$741 \quad \phi(i) \leq \phi(\psi(i)) \leq \phi(\psi^{(2)}(i)) \leq \dots \quad \forall i \in \llbracket 1, p \rrbracket \quad (22)$$

742 where the superscript denotes composition and where we can apply the permutation ψ an arbitrary number of times.

743 **\mathbf{Q} must be a sign matrix** Suppose that ψ is not the identity. We can pick an index $k \in \llbracket 1, p \rrbracket$ where $k \neq \psi(k)$. Because
 744 ϕ and ψ are injective, we can apply them to the inequality any number of times $n \in \mathbb{N}$ and get

$$745 \quad \phi(\psi^{(n)}(k)) \neq \phi(\psi^{(n+1)}(k)) \quad (23)$$

746 which together with Eq. 22 implies

$$747 \quad \phi(k) < \phi(\psi(k)) < \phi(\psi^{(2)}(k)) < \dots \quad (24)$$

748 which can be applied any number of times. However, each inequality increases the index by at least one, while the index
 749 cannot go above p . Thus, applying the chain at least p times, we have a contradiction. So, the permutation ψ in \mathbf{Q} must be
 750 the identity, and \mathbf{Q} boils down to a sign matrix, *i.e.* a diagonal matrix with 1 and -1 on the diagonal.
 751

Q must be the identity Since $Q^\top = Q$, Eq. 18 can be reformulated as

$$P'^\top L' P' = D' Q D^{-1} P^\top L P \quad (25)$$

where we now know that $D' Q D^{-1}$ is a diagonal matrix. The matrix $P'^\top L' P'$ has a diagonal of ones and so too does $P^\top L P$. It follows that $D' Q D^{-1} = I$. Since D and D' have positive diagonal entries, it follows that the entries of Q cannot equal -1 . So, $Q = I$, and thus $D' = D$. This concludes the proof of the Lemma.

A.3.2. PROOF OF LEMMA 6

From the equality $P^\top T P = P'^\top T' P'$, we deduce that

$$T = \tilde{P}^\top T' \tilde{P} \quad (26)$$

where $\tilde{P} = P' P^\top$ is the permutation matrix represented by function ϕ , which is defined such that: $\forall k \in \llbracket 1, p \rrbracket$, $\tilde{P}_{k, \phi(k)} = 1$ and $\forall l \neq \phi(k)$, $\tilde{P}_{k, l} = 0$. Using the notation ϕ instead of \tilde{P} , Eq. 26 can be rewritten as, $\forall k, l \in \llbracket 1, p \rrbracket$,

$$T_{k, l} = T'_{\phi(k), \phi(l)} \quad (27)$$

We proceed by a proof by contradiction: assume that $\tilde{P} \neq I$. As a consequence, there exists an index k such that $\phi(k) \neq k$. Let us fix such an index as k . We can assume without loss of generality that

$$\phi(k) > k \quad (28)$$

otherwise we just have to invert signs and switch rows and columns in the following. By assumption, the strictly lower triangular part of T only has non-zero elements, so we have

$$T_{\phi(k), k} \neq 0 \quad (29)$$

Yet, from Eq. 27, we know that $T_{\phi(k), k} = T'_{\phi(\phi(k)), \phi(k)}$ and all the non-zero elements of T' are in its strictly lower triangular part, which implies

$$\phi(\phi(k)) > \phi(k) \quad (30)$$

This logic can be applied repeatedly as in the preceding lemma's proof, and we see that

$$\phi(k) < \phi^{(2)}(k) < \phi^{(3)}(k) < \dots \quad (31)$$

Now, this leads to an infinite strictly increasing sequence, which is contradictory since the index cannot grow greater than p . So, $\tilde{P} = I$, which means that $P' P^\top P = P$, and thus, using the orthogonality of P ,

$$P' = P \quad (32)$$

It also follows that $T' = T$, which concludes the proof of Lemma 6.

825 B. Likelihood-based criterion

826 B.1. Basic setting

827 The fundamental question here is to choose the causal direction for two variables, having many views i . We denote the two
828 variables by x and y to avoid too many indices. It could be that the causal direction is $x \rightarrow y$:

$$830 y^i = b_i x^i + e^i \quad \text{for all } i \quad (33)$$

831 or $y \rightarrow x$:

$$832 x^i = c_i y^i + d^i \quad \text{for all } i \quad (34)$$

833 Recall that superscripts are for view in the case of the random variables; however, for the regression coefficients we use
834 subscripts because we need to take squares of such quantities. For the direction $x \rightarrow y$, collect the regressor and residuals
835 in the vectors $\mathbf{x} = (x^1, \dots, x^m)^\top$ and $\mathbf{e} = (e^1, \dots, e^m)^\top$, respectively; and likewise for \mathbf{y} and \mathbf{d} in the opposite direction.
836 Second-order moment matrices are denoted by $\Sigma_x = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$, $\Sigma_y = \mathbb{E}[\mathbf{y}\mathbf{y}^\top]$, and likewise for \mathbf{d} and \mathbf{e} . Note that the
837 independence between \mathbf{x} and \mathbf{e} and between \mathbf{y} and \mathbf{d} , and the fact that \mathbf{x} and \mathbf{y} are standardized, imply that

$$838 b_i = c_i = \text{cov}(x_i, y_i) \quad \text{for all } i \quad (35)$$

842 B.2. Derivation of the expected log-likelihood in one direction

843 We start by formulating the likelihood of the model in the case where disturbances are assumed to be Gaussian.

844 Consider the log-likelihood of (\mathbf{x}, \mathbf{y}) for the direction $x \rightarrow y$. Under the model, the prior (marginal) distribution of \mathbf{x} is
845 $\mathcal{N}(\mathbf{0}, \Sigma_x)$, and the prior distribution of $\mathbf{e} = \mathbf{y} - \mathbf{b} \odot \mathbf{x}$ is $\mathcal{N}(\mathbf{0}, \Sigma_e)$ — where \odot denotes the elementwise product. Thus,
846 we have

$$847 \log p(\mathbf{x}, \mathbf{y}; \mathbf{b}, \Sigma_x, \Sigma_e, x \rightarrow y) = \log p(\mathbf{y}|\mathbf{x}; \mathbf{b}, \Sigma_x, \Sigma_e, x \rightarrow y) + \log p(\mathbf{x}; \mathbf{b}, x \rightarrow y) \quad (36)$$

$$848 = -(\mathbf{y} - \mathbf{b} \odot \mathbf{x})^\top \Sigma_e^{-1} (\mathbf{y} - \mathbf{b} \odot \mathbf{x}) - \log |\Sigma_e| - \mathbf{x}^\top \Sigma_x^{-1} \mathbf{x} - \log |\Sigma_x| \quad (37)$$

$$849 = -\text{tr}(\Sigma_e^{-1} (\mathbf{y} - \mathbf{b} \odot \mathbf{x})(\mathbf{y} - \mathbf{b} \odot \mathbf{x})^\top) - \log |\Sigma_e| - \text{tr}(\Sigma_x^{-1} \mathbf{x}\mathbf{x}^\top) - \log |\Sigma_x| \quad (38)$$

850 where tr is the trace operator, and where we neglect the part of the normalization constant that does not depend on any
851 parameters.

852 The expected log-likelihood is thus

$$853 \mathbb{E}[\log p(\mathbf{x}, \mathbf{y}; \mathbf{b}, \Sigma_x, \Sigma_e, x \rightarrow y)] = -\text{tr}(\Sigma_e^{-1} \mathbb{E}[(\mathbf{y} - \mathbf{b} \odot \mathbf{x})(\mathbf{y} - \mathbf{b} \odot \mathbf{x})^\top]) - \log |\Sigma_e|$$

$$854 - \text{tr}(\Sigma_x^{-1} \mathbb{E}[\mathbf{x}\mathbf{x}^\top]) - \log |\Sigma_x| \quad (39)$$

$$855 = -\log |\Sigma_x| - \log |\Sigma_e| - 2m \quad (40)$$

856 Denote by \mathcal{L} this expected log-likelihood. We have (up to constants):

$$857 \mathcal{L}(\Sigma_x, \Sigma_e; x \rightarrow y) = -\log |\Sigma_x| - \log |\Sigma_e| \quad (41)$$

862 B.3. Derivation of the likelihood-based criterion

863 We now derive a criterion that can be used to find the causal direction.

864 The expected log-likelihood of the model for the direction $y \rightarrow x$ can be obtained by switching the roles of x and y
865 in Eq. 41. Thus, we have

$$866 \mathcal{L}(\Sigma_y, \Sigma_d; y \rightarrow x) = -\log |\Sigma_y| - \log |\Sigma_d| \quad (42)$$

867 where we recall that $\Sigma_y = \mathbb{E}[\mathbf{y}\mathbf{y}^\top]$ and $\Sigma_d = \mathbb{E}[\mathbf{d}\mathbf{d}^\top]$.

868 The correct direction can then be found by comparing the expected log-likelihoods of both directions. In particular, we
869 calculate the difference between the two log-likelihoods, which can be reformulated as a likelihood ratio (LR):

$$870 \text{LR} := \mathcal{L}(\Sigma_x, \Sigma_e; x \rightarrow y) - \mathcal{L}(\Sigma_y, \Sigma_d; y \rightarrow x) = \mathbb{E} \left[\log \frac{p(\mathbf{x}, \mathbf{y}; \mathbf{b}, \Sigma_x, \Sigma_e, x \rightarrow y)}{p(\mathbf{x}, \mathbf{y}; \mathbf{c}, \Sigma_y, \Sigma_d, y \rightarrow x)} \right] \quad (43)$$

Using Eq. 41 and Eq. 42, the criterion LR becomes

$$\text{LR} = -\log |\Sigma_x| - \log |\Sigma_e| + \log |\Sigma_y| + \log |\Sigma_d| . \quad (44)$$

Intuitively, the criterion in Eq. 44 will be large if the x^i are highly correlated (more than the y^i) over the views, since then the first determinant will be small. This should make intuitive sense since it means the cause is highly structured. Likewise, if the residuals for $x \rightarrow y$ are highly structured, that direction is more likely. In other words, the direction that exhibits more structure wins.

Surprisingly, in this formulation, there is actually no particular model for the dependencies of the disturbances: the likelihood should simply prefer dependent disturbances. It is important to note that the criterion makes sense even if the distribution is non-Gaussian, as will be proven below. We thus obtain a method that works for data of any distribution, but using only second-order statistics.

B.4. Non-negativity of the likelihood-based criterion

In this section, we prove that the criterion in Eq. 44 is non-negative for the direction $x \rightarrow y$. Importantly, no assumption of Gaussianity needs to be made here.

Consider the direction $x \rightarrow y$, and let us write $\mathbf{B} := \text{diag}(b_1, \dots, b_m)$. Using the matrix \mathbf{B} rather than the vector \mathbf{b} makes notations clearer in the proofs; please note that this notation is different from the main paper. In the following, we use interchangeably \mathbf{b} and \mathbf{B} . The covariance matrix of $\mathbf{y} = \mathbf{B}\mathbf{x} + \mathbf{e}$ is

$$\Sigma_y = \mathbf{B}\Sigma_x\mathbf{B} + \Sigma_e . \quad (45)$$

Given that $b_i = c_i$ for all i , we have $\mathbf{d} = \mathbf{x} - \mathbf{B}\mathbf{y} = \mathbf{x} - \mathbf{B}(\mathbf{B}\mathbf{x} + \mathbf{e}) = (\mathbf{I} - \mathbf{B}^2)\mathbf{x} - \mathbf{B}\mathbf{e}$. Thus, the covariance matrix of \mathbf{d} is

$$\Sigma_d = (\mathbf{I} - \mathbf{B}^2)\Sigma_x(\mathbf{I} - \mathbf{B}^2) + \mathbf{B}\Sigma_e\mathbf{B} . \quad (46)$$

Using these expressions, the criterion in Eq. 44 becomes

$$\text{LR} = -\log |\Sigma_x| - \log |\Sigma_e| + \log |\mathbf{B}\Sigma_x\mathbf{B} + \Sigma_e| + \log |(\mathbf{I} - \mathbf{B}^2)\Sigma_x(\mathbf{I} - \mathbf{B}^2) + \mathbf{B}\Sigma_e\mathbf{B}| . \quad (47)$$

Collect the variances of the residuals \mathbf{e} in the diagonal matrix

$$\mathbf{L} := \text{diag} \left(\sqrt{\text{var}(e^1)}, \dots, \sqrt{\text{var}(e^m)} \right) \quad (48)$$

and define the correlation matrix of the residuals as $\tilde{\Sigma}_e$, so that

$$\Sigma_e = \mathbf{L}\tilde{\Sigma}_e\mathbf{L} . \quad (49)$$

Now, the objective becomes

$$\begin{aligned} \text{LR} = & -\log |\Sigma_x| - \log |\tilde{\Sigma}_e| - \log |\mathbf{L}^2| + \log |\mathbf{B}\Sigma_x\mathbf{B} + \mathbf{L}\tilde{\Sigma}_e\mathbf{L}| \\ & + \log |(\mathbf{I} - \mathbf{B}^2)\Sigma_x(\mathbf{I} - \mathbf{B}^2) + \mathbf{B}\mathbf{L}\tilde{\Sigma}_e\mathbf{L}\mathbf{B}| . \end{aligned} \quad (50)$$

Importantly, the fact that x and y are standardized means that Σ_x and Σ_y have unit diagonal, so Eq. 45 implies

$$\mathbf{L}^2 + \mathbf{B}^2 = \mathbf{I} \quad (51)$$

and it follows that the diagonal entries of \mathbf{L} and \mathbf{B} are in $[-1, 1]$. Thus, we have

$$\text{LR} = -\log |\Sigma_x| - \log |\tilde{\Sigma}_e| - \log |\mathbf{L}^2| + \log |\mathbf{B}\Sigma_x\mathbf{B} + \mathbf{L}\tilde{\Sigma}_e\mathbf{L}| + \log |\mathbf{L}^2\Sigma_x\mathbf{L}^2 + \mathbf{L}\mathbf{B}\tilde{\Sigma}_e\mathbf{B}\mathbf{L}| \quad (52)$$

$$= -\log |\Sigma_x| - \log |\tilde{\Sigma}_e| + \log |\mathbf{B}\Sigma_x\mathbf{B} + \mathbf{L}\tilde{\Sigma}_e\mathbf{L}| + \log |\mathbf{L}\Sigma_x\mathbf{L} + \mathbf{B}\tilde{\Sigma}_e\mathbf{B}| \quad (53)$$

where \mathbf{L} and \mathbf{B} commuted in the first equality because they are diagonal.

The following lemma (proven in Appendix B.5) shows that this criterion is non-negative, and even positive under some assumptions.

Lemma 8 Let Σ_x and $\widetilde{\Sigma}_e$ be $m \times m$ real symmetric positive definite matrices with unit diagonal entries. Let $\mathbf{B} = \text{diag}(b_1, \dots, b_m)$ and $\mathbf{L} = \text{diag}(l_1, \dots, l_m)$ be diagonal matrices with entries in $[-1, 1]$, satisfying $\mathbf{B}^2 + \mathbf{L}^2 = \mathbf{I}$ (i.e. $b_i^2 + l_i^2 = 1$ for all i). Define the function

$$J(\mathbf{B}, \mathbf{L}) := -\log \det \Sigma_x - \log \det \widetilde{\Sigma}_e + \log \det (\mathbf{B}\Sigma_x\mathbf{B} + \mathbf{L}\widetilde{\Sigma}_e\mathbf{L}) + \log \det (\mathbf{L}\Sigma_x\mathbf{L} + \mathbf{B}\widetilde{\Sigma}_e\mathbf{B}). \quad (54)$$

Under the above assumptions, we have $J(\mathbf{B}, \mathbf{L}) \geq 0$. Moreover, $J(\mathbf{B}, \mathbf{L}) = 0$ if and only if

$$\mathbf{L}\widetilde{\Sigma}_e\mathbf{B} = \mathbf{B}\Sigma_x\mathbf{L}. \quad (55)$$

B.5. Proof of a useful lemma

Let us prove Lemma 8.

Step 1: define \mathbf{A} and \mathbf{C} — Define

$$\mathbf{A} := \mathbf{B}\Sigma_x\mathbf{B} + \mathbf{L}\widetilde{\Sigma}_e\mathbf{L} \quad \text{and} \quad \mathbf{C} := \mathbf{L}\Sigma_x\mathbf{L} + \mathbf{B}\widetilde{\Sigma}_e\mathbf{B}. \quad (56)$$

We have

$$J(\mathbf{B}, \mathbf{L}) := -\log \det \Sigma_x - \log \det \widetilde{\Sigma}_e + \log \det \mathbf{A} + \log \det \mathbf{C}. \quad (57)$$

Step 2: define \mathbf{M} — Define the $2m \times 2m$ block matrix

$$\mathbf{M} := \begin{pmatrix} \mathbf{B} & \mathbf{L} \\ -\mathbf{L} & \mathbf{B} \end{pmatrix}. \quad (58)$$

Since \mathbf{B} and \mathbf{L} are diagonal, they commute, hence $\mathbf{B}\mathbf{L} = \mathbf{L}\mathbf{B}$. Using $\mathbf{B}^2 + \mathbf{L}^2 = \mathbf{I}$, we compute

$$\mathbf{M}\mathbf{M}^\top = \begin{pmatrix} \mathbf{B}\mathbf{B} + \mathbf{L}\mathbf{L} & -\mathbf{B}\mathbf{L} + \mathbf{L}\mathbf{B} \\ -\mathbf{L}\mathbf{B} + \mathbf{B}\mathbf{L} & \mathbf{L}\mathbf{L} + \mathbf{B}\mathbf{B} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} = \mathbf{I}_{2m}. \quad (59)$$

Thus \mathbf{M} is orthonormal, so $\det \mathbf{M} = \pm 1$.

Step 3: define \mathcal{B} — Consider the block-diagonal matrix $\text{diag}(\Sigma_x, \widetilde{\Sigma}_e)$ and conjugate by \mathbf{M} :

$$\mathcal{B} := \mathbf{M} \begin{pmatrix} \Sigma_x & \mathbf{0} \\ \mathbf{0} & \widetilde{\Sigma}_e \end{pmatrix} \mathbf{M}^\top = \begin{pmatrix} \mathbf{A} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{C} \end{pmatrix}, \quad \mathbf{S} := \mathbf{L}\widetilde{\Sigma}_e\mathbf{B} - \mathbf{B}\Sigma_x\mathbf{L}. \quad (60)$$

Because \mathbf{M} is orthonormal and $\Sigma_x, \widetilde{\Sigma}_e \succ \mathbf{0}$, the matrix \mathcal{B} is symmetric positive definite (SPD). In particular, \mathbf{A} and \mathbf{C} are also SPD (as they are principal blocks of a SPD matrix) and thus \mathbf{A}^{-1} exists.

Step 4: Schur complement of \mathbf{A} in \mathcal{B} — For a symmetric block matrix $\begin{pmatrix} \mathbf{A} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{C} \end{pmatrix}$ with $\mathbf{A} \succ \mathbf{0}$, the Schur complement of \mathbf{A} is $\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}$, and the Schur's formula yields

$$\det \begin{pmatrix} \mathbf{A} & \mathbf{S} \\ \mathbf{S}^\top & \mathbf{C} \end{pmatrix} = \det(\mathbf{A}) \det(\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}). \quad (61)$$

Applying this to \mathcal{B} and using multiplicativity of determinant and the fact that $\det \mathbf{M} = \pm 1$ gives

$$\det(\Sigma_x) \det(\widetilde{\Sigma}_e) = \det(\mathcal{B}) = \det(\mathbf{A}) \det(\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}). \quad (62)$$

Step 5: derive a matrix inequality — Because $\mathcal{B} \succ \mathbf{0}$ and $\mathbf{A} \succ \mathbf{0}$, it follows from the Schur complement characterization of positive definiteness that $\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S} \succ \mathbf{0}$. In particular, we have $\det(\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}) > 0$. Moreover, since $\mathbf{A} \succ \mathbf{0}$, we have $\mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S} \succeq \mathbf{0}$, and thus $\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S} \preceq \mathbf{C}$.

990 **Step 6: derive a determinant inequality** — If $\mathbf{0} \preceq \mathbf{X} \preceq \mathbf{Y}$ and $\mathbf{Y} \succ \mathbf{0}$, then conjugating by $\mathbf{Y}^{-1/2}$
 991 shows $\mathbf{0} \preceq \mathbf{Y}^{-1/2} \mathbf{X} \mathbf{Y}^{-1/2} \preceq \mathbf{I}$. Thus, all eigenvalues of $\mathbf{Y}^{-1/2} \mathbf{X} \mathbf{Y}^{-1/2}$ lie in $[0, 1]$, so their product satisfies
 992 $\det(\mathbf{Y}^{-1/2} \mathbf{X} \mathbf{Y}^{-1/2}) \leq 1$. Hence $\det(\mathbf{X}) \leq \det(\mathbf{Y})$. Applying this to $\mathbf{X} := \mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}$ and $\mathbf{Y} := \mathbf{C}$ yields

$$993 \det(\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}) \leq \det(\mathbf{C}) . \quad (63)$$

994
 995
 996 **Step 7: $J(\mathbf{B}, \mathbf{L}) \geq 0$** — From Steps 4 and 6, we get

$$997 \det(\boldsymbol{\Sigma}_x) \det(\widetilde{\boldsymbol{\Sigma}}_e) = \det(\mathbf{A}) \det(\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}) \leq \det(\mathbf{A}) \det(\mathbf{C}) . \quad (64)$$

998
 999 Taking natural logarithms and rearranging gives

$$1000 \log \det(\mathbf{A}) + \log \det(\mathbf{C}) \geq \log \det(\boldsymbol{\Sigma}_x) + \log \det(\widetilde{\boldsymbol{\Sigma}}_e) , \quad (65)$$

1001 which is precisely $J(\mathbf{B}, \mathbf{L}) \geq 0$.

1002
 1003
 1004 **Step 8: equality condition** — The equality $J(\mathbf{B}, \mathbf{L}) = 0$ holds iff

$$1005 \det(\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}) = \det(\mathbf{C}) . \quad (66)$$

1006 Write $\mathbf{K} := \mathbf{C}^{-1/2} (\mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}) \mathbf{C}^{-1/2} \succeq \mathbf{0}$. We have

$$1007 \det(\mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}) = \det(\mathbf{C}) \det(\mathbf{I} - \mathbf{C}^{-1} \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S}) \quad (67)$$

$$1008 = \det(\mathbf{C}) \det(\mathbf{C}^{-\frac{1}{2}} \mathbf{C}^{\frac{1}{2}} - \mathbf{C}^{-\frac{1}{2}} \mathbf{C}^{-\frac{1}{2}} \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S} \mathbf{C}^{-\frac{1}{2}} \mathbf{C}^{\frac{1}{2}}) \quad (68)$$

$$1009 = \det(\mathbf{C}) \det(\mathbf{C}^{-\frac{1}{2}}) \det(\mathbf{I} - \mathbf{K}) \det(\mathbf{C}^{\frac{1}{2}}) \quad (69)$$

$$1010 = \det(\mathbf{C}) \det(\mathbf{I} - \mathbf{K}) . \quad (70)$$

1011 Thus, the equality condition is $\det(\mathbf{I} - \mathbf{K}) = 1$. From Step 5, we know that $\mathbf{0} \preceq \mathbf{C} - \mathbf{S}^\top \mathbf{A}^{-1} \mathbf{S} \preceq \mathbf{C}$. Conjugating by
 1012 $\mathbf{C}^{-\frac{1}{2}}$ gives $\mathbf{0} \preceq \mathbf{I} - \mathbf{K} \preceq \mathbf{I}$, so each eigenvalue λ_i of \mathbf{K} satisfies $0 \leq \lambda_i \leq 1$. The eigenvalues of $\mathbf{I} - \mathbf{K}$ are $1 - \lambda_i \in [0, 1]$,
 1013 hence

$$1014 \prod_i (1 - \lambda_i) = \det(\mathbf{I} - \mathbf{K}) = 1 . \quad (71)$$

1015 The equality thus holds iff $\lambda_i = 0$ for all i , hence $\mathbf{K} = \mathbf{0}$. Since $\mathbf{C}^{-\frac{1}{2}}, \mathbf{A}^{-1} \succ \mathbf{0}$, we get $\mathbf{S} = \mathbf{0}$. Conversely, $\mathbf{S} = \mathbf{0}$
 1016 clearly forces equality. Therefore, the equality occurs exactly when

$$1017 \mathbf{S} = \mathbf{L} \widetilde{\boldsymbol{\Sigma}}_e \mathbf{B} - \mathbf{B} \boldsymbol{\Sigma}_x \mathbf{L} = \mathbf{0} , \quad (72)$$

1018 *i.e.*

$$1019 \forall i, j : \quad b_i l_j (\boldsymbol{\Sigma}_x)_{ij} = l_i b_j (\widetilde{\boldsymbol{\Sigma}}_e)_{ij} . \quad (73)$$

1020 This completes the proof.

1021 B.6. Conditions for strict positivity of the likelihood-based criterion

1022 From now on, we assume that the model is not degenerate in the sense that disturbances have a positive variance, that is
 1023 $l_i = \text{var}(e_i) > 0$ for all i .

1024 The criterion LR derived in Appendix B.4 is used to determine the causal direction from its sign, with LR = 0 indicating
 1025 that the two directions cannot be distinguished. Lemma 8 shows that LR = 0 if and only if

$$1026 \mathbf{L} \widetilde{\boldsymbol{\Sigma}}_e \mathbf{B} = \mathbf{B} \boldsymbol{\Sigma}_x \mathbf{L} . \quad (74)$$

1027 Entrywise, this is equivalent to

$$1028 \forall i, j : \quad l_i b_j (\widetilde{\boldsymbol{\Sigma}}_e)_{ij} = b_i l_j (\boldsymbol{\Sigma}_x)_{ij} . \quad (75)$$

1029 Let us fix $i, j \in [1, m]$.

1045 **Case 1:** $b_i = 0$ and $b_j = 0$ — The equality immediately holds.

1046
1047 **Case 2:** $b_i = 0$ and $b_j \neq 0$ — The right-hand side in Eq. 75 is equal to 0. Since $l_i = \text{var}(e^i) > 0$ and $b_j \neq 0$, we must
1048 have $(\widetilde{\Sigma}_e)_{ij} = 0$.

1049
1050 **Case 3:** $b_i \neq 0$ and $b_j = 0$ — The left-hand side in Eq. 75 is equal to 0. Since $l_j = \text{var}(e^j) > 0$ and $b_i \neq 0$, we must
1051 have $(\Sigma_x)_{ij} = 0$.

1052
1053 **Case 4:** $b_i \neq 0$ and $b_j \neq 0$ — Multiplying Eq. 75 by $\frac{l_j}{b_i b_j^2}$, and using $l_j^2 = 1 - b_j^2$, we obtain

$$\frac{l_i l_j}{b_i b_j} (\widetilde{\Sigma}_e)_{ij} = \left(\frac{1}{b_j^2} - 1 \right) (\Sigma_x)_{ij} . \quad (76)$$

1054
1055
1056 Switching the roles of i and j , and using the symmetry of $(\widetilde{\Sigma}_e)_{ij}$ and $(\Sigma_x)_{ij}$, yield

$$\frac{l_i l_j}{b_i b_j} (\widetilde{\Sigma}_e)_{ij} = \left(\frac{1}{b_i^2} - 1 \right) (\Sigma_x)_{ij} . \quad (77)$$

1057
1058
1059 hence

$$b_i^2 (\Sigma_x)_{ij} = b_j^2 (\Sigma_x)_{ij} . \quad (78)$$

1060
1061 Similarly, exchanging the roles of $(b_i, b_j, (\Sigma_x)_{ij})$ and $(l_i, l_j, (\widetilde{\Sigma}_e)_{ij})$ yields

$$l_i^2 (\widetilde{\Sigma}_e)_{ij} = l_j^2 (\widetilde{\Sigma}_e)_{ij} . \quad (79)$$

1062
1063
1064 Now suppose $(\Sigma_x)_{ij} \neq 0$ or $(\widetilde{\Sigma}_e)_{ij} \neq 0$. From Eq. 78 and Eq. 79 it follows that $b_i^2 = b_j^2$ or $l_i^2 = l_j^2$. Since $b_i^2 + l_i^2 = 1$,
1065 both equalities hold simultaneously:

$$|b_i| = |b_j| \quad \text{and} \quad l_i = l_j \quad (80)$$

1066
1067 where the absolute value can be dropped for l_i because $l_i = \text{var}(e^i) > 0$. Plugging these equalities into Eq. 76 yields

$$l_j^2 (\widetilde{\Sigma}_e)_{ij} = \pm (1 - b_j^2) (\Sigma_x)_{ij} \quad (81)$$

1068
1069 hence

$$|(\widetilde{\Sigma}_e)_{ij}| = |(\Sigma_x)_{ij}| . \quad (82)$$

1070
1071 Moreover, for Eq. 75 to hold we must also have the sign equality

$$\text{sign}(b_i) \text{sign}(b_j) = \text{sign}((\Sigma_x)_{ij}) \text{sign}((\widetilde{\Sigma}_e)_{ij}) . \quad (83)$$

1072
1073 Altogether, either both $(\Sigma_x)_{ij}$ and $(\widetilde{\Sigma}_e)_{ij}$ vanish, or the following four conditions hold:

$$|(\widetilde{\Sigma}_e)_{ij}| = |(\Sigma_x)_{ij}| \neq 0 \quad (84)$$

$$|b_i| = |b_j| \quad (85)$$

$$l_i = l_j \quad (86)$$

$$\text{sign}(b_i) \text{sign}(b_j) = \text{sign}((\Sigma_x)_{ij}) \text{sign}((\widetilde{\Sigma}_e)_{ij}) . \quad (87)$$

1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099 Conversely, if $(\Sigma_x)_{ij} = (\widetilde{\Sigma}_e)_{ij} = 0$, or if conditions Eq. 84–Eq. 87 are satisfied, then Eq. 75 holds.

Conclusion — By contraposition, Eq. 74 is violated whenever there exists (i, j) such that one of the following cases holds:

- $b_i = 0, b_j \neq 0$, and $(\widetilde{\Sigma}_e)_{ij} \neq 0$
- $b_i \neq 0, b_j = 0$, and $(\Sigma_x)_{ij} \neq 0$
- $b_i \neq 0, b_j \neq 0, (\Sigma_x)_{ij} \neq 0$ or $(\widetilde{\Sigma}_e)_{ij} \neq 0$, and at least one of the conditions Eq. 84–Eq. 87 is not met.

In particular, a simple sufficient condition is

$$\exists i, j : \quad (b_i \neq 0 \quad \text{or} \quad b_j \neq 0) \quad \text{and} \quad |(\widetilde{\Sigma}_e)_{ij}| \neq |(\Sigma_x)_{ij}| . \quad (88)$$

which is the basis of Assumptions 1 and 2 given in the main paper.

C. Cross-covariance-based criterion

Consider the setting described in Appendix B.1.

C.1. Derivation of the cross-covariance-based criterion

A simple approach to infer the causal direction is to exploit the assumption that root variables and residuals are independent. Specifically, if the true direction is $x \rightarrow y$, then \mathbf{x} and \mathbf{e} are independent; if the true direction is $y \rightarrow x$, then \mathbf{y} and \mathbf{d} are independent. This implies that the cross-moment matrix $\mathbb{E}[\mathbf{e}\mathbf{x}^\top]$ should vanish when $x \rightarrow y$, while $\mathbb{E}[\mathbf{d}\mathbf{y}^\top]$ should vanish when $y \rightarrow x$. Comparing the Frobenius norms of these matrices therefore provides a natural criterion, based solely on cross-covariances:

$$\text{FC} := \left\| \mathbb{E}[\mathbf{d}\mathbf{y}^\top] \right\|_F - \left\| \mathbb{E}[\mathbf{e}\mathbf{x}^\top] \right\|_F . \quad (89)$$

Intuitively, if this criterion is positive we deduce that $x \rightarrow y$, and if negative we deduce that $y \rightarrow x$.

C.2. Conditions for strict positivity of the cross-covariance criterion

Assume without loss of generality that the true causal direction is $x \rightarrow y$, and let us analyze the sign of the criterion FC.

Given that $x \rightarrow y$, the true model is

$$\mathbf{y} = \mathbf{B}\mathbf{x} + \mathbf{e} \quad (90)$$

where $\mathbf{B} = \text{diag}(\text{cov}(x^1, y^1), \dots, \text{cov}(x^m, y^m))$, and \mathbf{e} is independent from \mathbf{x} . So, we have

$$\mathbb{E}[\mathbf{x}\mathbf{e}^\top] = \mathbf{0} . \quad (91)$$

In the wrong direction, the residuals take the form

$$\mathbf{d} = \mathbf{x} - \mathbf{C}\mathbf{y}, \quad \mathbf{C} = \mathbf{B} . \quad (92)$$

The cross-covariance between observations \mathbf{y} and residuals \mathbf{d} is then

$$\mathbb{E}[\mathbf{y}\mathbf{d}^\top] = \mathbb{E}[\mathbf{y}(\mathbf{x} - \mathbf{B}\mathbf{y})^\top] \quad (93)$$

$$= \mathbb{E}[(\mathbf{B}\mathbf{x} + \mathbf{e})(\mathbf{x} - \mathbf{B}(\mathbf{B}\mathbf{x} + \mathbf{e}))^\top] \quad (94)$$

$$= \mathbb{E}[\mathbf{B}\mathbf{x}\mathbf{x}^\top(\mathbf{I} - \mathbf{B}^2)] + \mathbb{E}[\mathbf{e}\mathbf{x}^\top(\mathbf{I} - \mathbf{B}^2)] - \mathbb{E}[\mathbf{B}\mathbf{x}\mathbf{e}^\top\mathbf{B}] - \mathbb{E}[\mathbf{e}\mathbf{e}^\top\mathbf{B}] \quad (95)$$

$$= \mathbf{B}\Sigma_x(\mathbf{I} - \mathbf{B}^2) - \Sigma_e\mathbf{B} \quad (96)$$

where $\Sigma_x = \mathbb{E}[\mathbf{x}\mathbf{x}^\top]$ and $\Sigma_e = \mathbb{E}[\mathbf{e}\mathbf{e}^\top]$.

From Eq. 91 and the non-negativity of the Frobenius norm, we see that FC is always non-negative. Moreover, using Eq. 96, we deduce that FC = 0 if and only if

$$\mathbf{B}\Sigma_x(\mathbf{I} - \mathbf{B}^2) - \Sigma_e\mathbf{B} = \mathbf{0} . \quad (97)$$

Note that the fact that $l_i > 0$, for all i , makes \mathbf{L} invertible. So, using the identity $\mathbf{I} - \mathbf{B}^2 = \mathbf{L}^2$, and multiplying on the right by \mathbf{L}^{-1} , we obtain

$$\mathbf{B}\Sigma_x\mathbf{L} = \Sigma_e\mathbf{B}\mathbf{L}^{-1} \quad (98)$$

hence

$$\mathbf{B}\Sigma_x\mathbf{L} = \mathbf{L}\tilde{\Sigma}_e\mathbf{B} \quad (99)$$

since $\Sigma_e = \mathbf{L}\tilde{\Sigma}_e\mathbf{L}$, and \mathbf{B} and \mathbf{L} commute.

This condition coincides exactly with Eq. 74. In particular, this shows that the likelihood-based criterion and the cross-covariance-based criterion share the same consistency assumptions.

D. Identifiability of LiMVAM

D.1. Multivariate conditions for strict positivity of a criterion

The following assumption is a multivariate generalization of the assumption detailed in B.6 and C.2.

Assumption 6 (Diversity of the views) *Let $j \neq j'$ denote any two component indices such that $\mathbf{x}_j \rightarrow \mathbf{x}_{j'}$. Then, there exist two views i and i' such that the three following conditions hold:*

1. $\text{corr}(\mathbf{x}_j^i, \mathbf{x}_{j'}^i) \neq 0$ or $\text{corr}(\mathbf{x}_j^{i'}, \mathbf{x}_{j'}^{i'}) \neq 0$
2. $\text{corr}(\mathbf{x}_j^i, \mathbf{x}_{j'}^{i'}) \neq 0$ or $\text{corr}(\mathbf{x}_j^{i'}, \mathbf{x}_{j'}^i) \neq 0$
3. At least one of the four following inequalities is met:

$$\begin{aligned} |\text{corr}(\mathbf{x}_j^i, \mathbf{x}_{j'}^{i'})| &\neq |\text{corr}(\mathbf{x}_j^{i'}, \mathbf{x}_{j'}^i)|, |\text{corr}(\mathbf{x}_j^i, \mathbf{x}_{j'}^i)| \neq |\text{corr}(\mathbf{x}_j^{i'}, \mathbf{x}_{j'}^{i'})|, \text{var}(\mathbf{e}_{j'}^i) \neq \text{var}(\mathbf{e}_{j'}^{i'}) \\ \text{sign}(\text{corr}(\mathbf{x}_j^i, \mathbf{x}_{j'}^{i'})) \text{sign}(\text{corr}(\mathbf{x}_j^{i'}, \mathbf{x}_{j'}^i)) &\neq \text{sign}(\text{corr}(\mathbf{x}_j^i, \mathbf{x}_{j'}^i)) \text{sign}(\text{corr}(\mathbf{x}_j^{i'}, \mathbf{x}_{j'}^{i'})) \end{aligned}$$

where corr denotes the correlation coefficient and sign denotes the sign function.

These conditions can only be met for indices $i \neq i'$, so they can be interpreted as a need for correlation and diversity across views. Because they may be hard to interpret at first glance, we next visualize in Figure 5 what these conditions mean for the causal graph.

In particular, Assumption 6 immediately implies its simpler version used in the main text, Assumption 2.

D.2. An example for interpreting the SOS-based identifiability condition in Assumption 2

Consider one view and two variables denoted for notational simplicity x and y . We have

$$y^1 = b^1 x^1 + e^1 \quad (100)$$

and assume for simplicity that x^1 and y^1 have been standardized.

Now, suppose we observe two datasets from this one view, denote them by (x^1, y^1) and $(x^{1'}, y^{1'})$, which have identical distributions and are independent of each other. Crucially, all this data is from one view. However, we could artificially generate a new “view” by adding those two datasets together:

$$x^2 = x^1 + x^{1'} \quad (101)$$

and likewise for y^2 and e^2 . Clearly the new data follows the SEM with the same b :

$$y^2 = b^1 x^2 + e^2 \quad (102)$$

since we are simply adding each of the terms separately in the two datasets.

But this “multi-view” data is degenerate. In fact, the correlation coefficients are

$$\text{corr}(x^1, x^2) = \frac{\text{cov}(x^1, x^1) + \text{cov}(x^1, x^{1'})}{\text{std}(x^1)\text{std}(x^2)} = \frac{\text{var}(x^1)}{\sqrt{2}\text{std}(x^1)\text{std}(x^1)} = \frac{1}{\sqrt{2}}. \quad (103)$$

An exactly identical calculation applies for $\text{corr}(e^1, e^2)$ which is also equal to $1/\sqrt{2}$.

Thus, we see that this case violates the simple identifiability condition in Assumption 2. This is understandable since no new information is created when computing the variable in the second view, and the two views are thus degenerate.

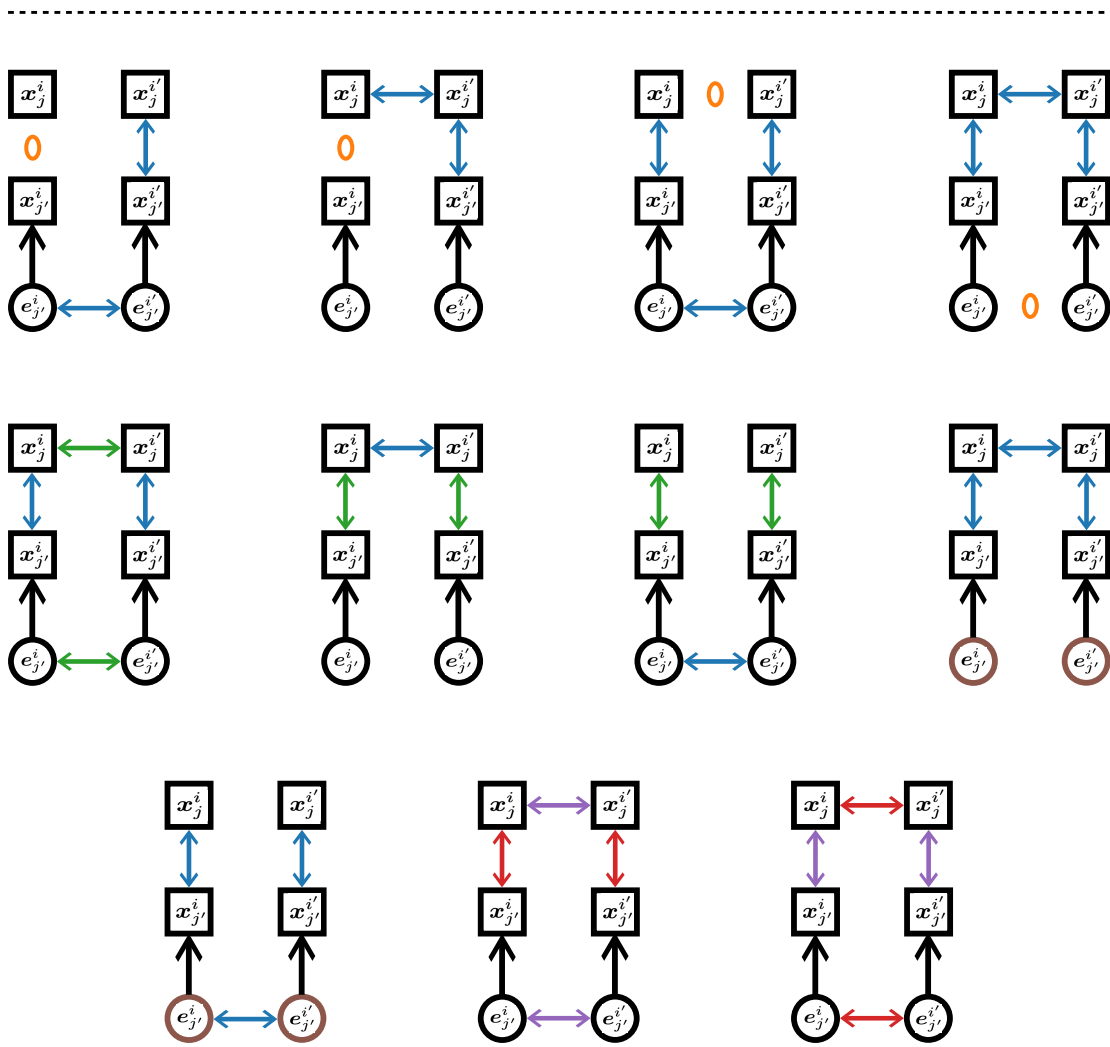
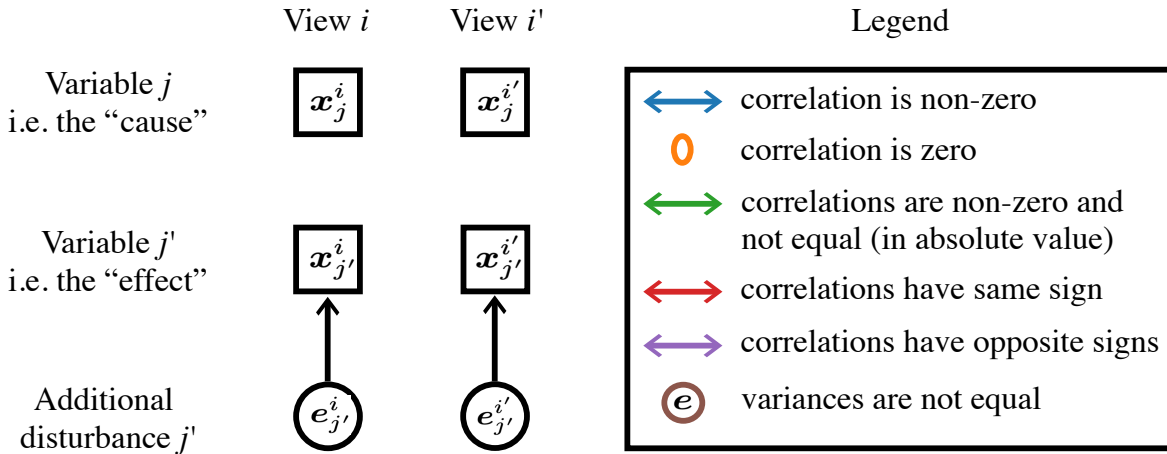


Figure 5. Assume the direction is $x_j \rightarrow x_{j'}$. Assumption 6 corresponds to being in one of the 11 situations showed in the figure. Note that no arrow means that no information is required.

D.3. Proof of Theorem 3

We proceed in three steps. (i) We show that the LiMVAM model always admits at least one root variable and that our LR and FC criteria recover one consistently. (ii) We prove a *recursive residuals* lemma: removing a recovered root and regressing on it preserves the model form and the causal ordering, which implies that the full (partial) ordering is consistently recoverable. (iii) Conditional on the recovered ordering, each view-specific strictly lower-triangular system can be estimated consistently, hence the adjacency matrices are identifiable.

Preliminaries For clarity, we unpack the “shared ordering” assumption. This assumption is equivalent to the existence of a *main* DAG \mathcal{G} that is the union of view-specific DAGs $\{\mathcal{G}^i\}_{i=1}^m$: an edge appears in \mathcal{G} if it appears in at least one \mathcal{G}^i . Root variables are therefore common candidates across views. Furthermore, in terms of vocabulary, we distinguish (a) *directed paths* (a sequence of edges with consistent direction) from (b) *direct edges* (one edge between adjacent variables). An *indirect path* is a directed path with at least two edges.

D.3.1. RECOVERING A ROOT VARIABLE

Let observations follow LiMVAM in Eq. 3. Since all views share an ordering, there exists a (possibly non-unique) permutation \mathbf{P} such that

$$\mathbf{T}^i = \mathbf{P}^\top \mathbf{B}^i \mathbf{P} \quad \text{is strictly lower triangular for all } i \in \llbracket 1, m \rrbracket . \quad (104)$$

Reordering the observations and disturbances with $\mathbf{x}^{\prime i} = \mathbf{P}\mathbf{x}^i$ and $\mathbf{e}^{\prime i} = \mathbf{P}\mathbf{e}^i$, the model is

$$\mathbf{x}^{\prime i} = \mathbf{T}^i \mathbf{x}^{\prime i} + \mathbf{e}^{\prime i}, \quad i \in \llbracket 1, m \rrbracket . \quad (105)$$

The first row of each \mathbf{T}^i is zero, so $x_1^{\prime i}$ is exogenous; hence a root variable always exists.

To *find* a root, consider distinct indices (j, k) and define M_{jk} to be the criterion from Eq. 5 or Eq. 6 (with $M_{kj} = -M_{jk}$). Stack these into $\mathbf{M} \in \mathbb{R}^{p \times p}$ with zero diagonal, and analyze signs under Assumption 6 (which concerns only direct relations):

1. **Direct edge:** $x_j \rightarrow x_k$ or $x_k \rightarrow x_j$. Under the bivariate conditions derived from Assumption 6, $M_{jk} > 0 > M_{kj}$ when $x_j \rightarrow x_k$, and $M_{jk} < 0 < M_{kj}$ when $x_k \rightarrow x_j$.
2. **Indirect path:** $x_j \rightarrow \dots \rightarrow x_k$ or $x_k \rightarrow \dots \rightarrow x_j$. Consider the reduced form of the model in Eq. 4, $\mathbf{x}^i = \mathbf{A}^i \mathbf{e}^i$ with $\mathbf{A}^i = (\mathbf{I} - \mathbf{B}^i)^{-1}$. The entry A_{kj}^i is the *total causal effect* of x_j^i on x_k^i , i.e. the sum over all directed paths from j to k of the products of edge coefficients. Hence for each view i ,

$$x_k^i = A_{kj}^i x_j^i + \epsilon^i \quad (106)$$

where ϵ^i collects all terms not caused by x_j^i and is independent of x_j^i . Thus, the pair (x_j^i, x_k^i) behaves as if there were a direct edge with coefficient A_{kj}^i , and the sign analysis of Item 1 applies but with non-strict inequalities: $M_{jk} \geq 0 \geq M_{kj}$ when $x_j \rightarrow \dots \rightarrow x_k$ (and conversely). If the bivariate conditions hold for (j, k) , the inequalities are strict; however, we do not require this here.

3. **No relation, no common ancestor:** then x_j and x_k depend on disjoint disturbance sets and are independent, hence $M_{jk} = M_{kj} = 0$.
4. **No relation but with a common ancestor.** If x_j and x_k are not connected by any directed path but share at least one ancestor, then their pairwise criteria M_{jk} and M_{kj} may take mixed signs. However, since each of x_j and x_k has at least one parent, there exist variables $x_{j'}$ and $x_{k'}$ such that $M_{jj'} < 0$ and $M_{kk'} < 0$. Thus, even if the signs of M_{jk} and M_{kj} are not determined, the j -th and k -th rows of \mathbf{M} necessarily contain negative entries. For instance, with three variables, $x_1 \rightarrow x_2$ and $x_1 \rightarrow x_3$ makes x_2 and x_3 siblings: there is no directed path between them, but both rows 2 and 3 of \mathbf{M} contain negative entries due to their direct parent x_1 .

Consequently, x_j is a root (no incoming edges) *if and only if* the j -th row of \mathbf{M} has only nonnegative entries. Hence a root exists and is consistently detectable. If multiple roots exist, we pick one at random.

D.3.2. RECOVERING THE CAUSAL ORDERING

We now show the recursive step.

Lemma 9 (Residuals preserve LiMVAM and the ordering) *Let j be a root variable. Regress every other variable on x_j (within each view) and remove x_j . The residual vector still follows a LiMVAM model with the same ordering over the remaining variables.*

Proof. The proof is a direct adaptation of Shimizu et al. (2011, Lemma 2 and Corollary 1) to the LiMVAM model: the original result is in the single-view setting with non-Gaussian disturbances.

Choose a permutation P that places x_j first for simplicity (since x_j is a root variable, this permutation exists). In the permuted system we write

$$\mathbf{x}'^i = \mathbf{T}^i \mathbf{x}'^i + \mathbf{e}'^i, \quad \mathbf{x}^i = P \mathbf{x}'^i, \quad \mathbf{e}^i = P \mathbf{e}'^i \quad (107)$$

where each \mathbf{T}^i is strictly lower triangular and $\mathbf{x}'_1 = x_j$. Equivalently, the model can be expressed in reduced form

$$\mathbf{x}'^i = \mathbf{A}^i \mathbf{e}'^i, \quad \mathbf{A}^i := (\mathbf{I} - \mathbf{T}^i)^{-1} \quad (108)$$

where the matrix \mathbf{A}^i is lower triangular with unit diagonal.

The entry A_{k1}^i ($k \geq 2$) quantifies the dependence of x_k^i on e_1^i . Since $x'_1 = e_1^i$ is exogenous, we have

$$\mathbb{E}[x_k^i x_1^i] = A_{k1}^i \mathbb{E}[(x_1^i)^2] \quad (109)$$

so the least-squares regression coefficient of x_k^i on x_1^i equals A_{k1}^i . Therefore, regressing each x_k^i ($k \geq 2$) on x_1^i yields residuals $\mathbf{r}_{(1)}^i \in \mathbb{R}^{p-1}$ such that

$$\mathbf{r}_{(1)}^i = \mathbf{x}_{2:p}^i - \mathbf{A}_{2:p,1}^i x_1^i = \mathbf{A}_{2:p,:}^i \mathbf{e}'^i - \mathbf{A}_{2:p,1}^i e_1^i = \mathbf{A}_{2:p,2:p}^i \mathbf{e}'_{2:p} \quad (110)$$

Here $\mathbf{A}_{2:p,2:p}^i$ is again lower triangular with unit diagonal, and $\mathbf{e}'_{2:p}$ has mutually independent components. Hence $\mathbf{r}_{(1)}^i$ follows a LiMVAM model on $p-1$ variables. Moreover, since $\mathbf{A}_{2:p,2:p}^i$ is precisely the submatrix of \mathbf{A}^i obtained by deleting the first row and column, the relative structure among the remaining variables is unchanged. In other words, the causal ordering of variables $2, \dots, p$ is preserved. \square

Applying Lemma 9 recursively — find a root via M (under Assumption 6), regress out its effect, remove it, and repeat — recovers the full (partial) causal ordering consistently.

D.3.3. RECOVERING THE ADJACENCY MATRICES

Sections D.3.1 and D.3.2 showed that, under Assumption 6, the algorithms PairwiseLiMVAM and DirectLiMVAM consistently recover the causal ordering. Equivalently, they recover a permutation matrix P that makes the $\mathbf{T}^i = P^\top \mathbf{B}^i P$ strictly lower triangular. So, we can reorder the variables as follows:

$$\mathbf{x}'^i = \mathbf{T}^i \mathbf{x}'^i + \mathbf{e}'^i, \quad i \in \llbracket 1, m \rrbracket \quad (111)$$

Fix $j \in \llbracket 1, p \rrbracket$. We have

$$x_j^i = \sum_{k=1}^{j-1} T_{jk}^i x_k^i + e_j^i, \quad i \in \llbracket 1, m \rrbracket \quad (112)$$

Because e_j^i is independent of $(x_1^i, \dots, x_{j-1}^i)$, the regressors are exogenous. This is a system of linear regressions across views with potentially cross-view-correlated errors. Feasible GLS is consistent in this setting: once a consistent estimate of the cross-view error covariance is obtained, the GLS estimator converges to the true coefficients (Greene, 2003). Therefore, $\{T_{jk}^i\}_{k < j}$ are consistently estimated for every (i, j) , hence each \mathbf{T}^i is consistently recovered. Finally, $\mathbf{B}^i = P \mathbf{T}^i P^\top$, so the view-specific adjacency matrices are identifiable. This concludes the proof.

E. SOS-based Algorithms for LiMVAM

E.1. Consistent estimator of the likelihood-based criterion

Assume we have observed n independent samples, concatenated as $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ and $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]$, with sample index j .

In practice, the four covariance matrices in Eq. 44 can be estimated by consistent estimators from the samples in \mathbf{X} and \mathbf{Y} . Since $b_i = c_i = \text{cov}(x_i, y_i)$ for all i , we define

$$\forall i : \quad \hat{b}_i := \hat{c}_i := \frac{1}{n} \sum_{j=1}^n x_{ij} y_{ij} \xrightarrow{p} \text{cov}(x_i, y_i) = b_i = c_i \quad (113)$$

so that $\hat{\mathbf{b}} := (\hat{b}_1, \dots, \hat{b}_m)^\top \xrightarrow{p} \mathbf{b}$ and $\hat{\mathbf{c}} := (\hat{c}_1, \dots, \hat{c}_m)^\top \xrightarrow{p} \mathbf{c}$. Using these vectors, we define the empirical covariances of the residuals

$$\widehat{\Sigma}_e := \frac{1}{n} \sum_{j=1}^n (\mathbf{y}_j - \hat{\mathbf{b}} \odot \mathbf{x}_j)(\mathbf{y}_j - \hat{\mathbf{b}} \odot \mathbf{x}_j)^\top \xrightarrow{p} \Sigma_e \quad (114)$$

$$\widehat{\Sigma}_d := \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \hat{\mathbf{b}} \odot \mathbf{y}_j)(\mathbf{x}_j - \hat{\mathbf{b}} \odot \mathbf{y}_j)^\top \xrightarrow{p} \Sigma_d \quad (115)$$

and we also define the empirical covariances of the observations

$$\widehat{\Sigma}_x := \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j^\top \xrightarrow{p} \Sigma_x \quad (116)$$

$$\widehat{\Sigma}_y := \frac{1}{n} \sum_{j=1}^n \mathbf{y}_j \mathbf{y}_j^\top \xrightarrow{p} \Sigma_y \quad (117)$$

We can then evaluate the criterion in these estimators.

$$\widehat{\text{LR}} := \mathcal{L}(\widehat{\Sigma}_x, \widehat{\Sigma}_e; x \rightarrow y) - \mathcal{L}(\widehat{\Sigma}_y, \widehat{\Sigma}_d; y \rightarrow x) = -\log |\widehat{\Sigma}_x| - \log |\widehat{\Sigma}_e| + \log |\widehat{\Sigma}_y| + \log |\widehat{\Sigma}_d| \quad (118)$$

and it follows that this empirical criterion is a consistent estimator of the true criterion:

$$\mathcal{L}(\widehat{\Sigma}_x, \widehat{\Sigma}_e; x \rightarrow y) - \mathcal{L}(\widehat{\Sigma}_y, \widehat{\Sigma}_d; y \rightarrow x) \xrightarrow{p} \mathcal{L}(\Sigma_x, \Sigma_e; x \rightarrow y) - \mathcal{L}(\Sigma_y, \Sigma_d; y \rightarrow x) =: \text{LR} \quad (119)$$

E.2. Consistent estimator of the cross-covariance-based criterion

Assume we have observed n independent samples, concatenated as $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ and $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_n]$, with sample index j .

In practice, the matrices $\mathbb{E}[\mathbf{e}\mathbf{x}^\top]$ and $\mathbb{E}[\mathbf{d}\mathbf{y}^\top]$ are estimated from the data. Performing univariate least-squares regression of \mathbf{y} on \mathbf{x} (resp. \mathbf{x} on \mathbf{y}) yields the residual matrices $\widehat{\mathbf{E}} = [\widehat{\mathbf{e}}_1, \dots, \widehat{\mathbf{e}}_n]$ (resp. $\widehat{\mathbf{D}} = [\widehat{\mathbf{d}}_1, \dots, \widehat{\mathbf{d}}_n]$). The corresponding cross-covariance estimators are

$$\frac{1}{n} \sum_{j=1}^n \widehat{\mathbf{e}}_j \mathbf{x}_j^\top \xrightarrow{p} \mathbb{E}[\mathbf{e}\mathbf{x}^\top] \quad \text{and} \quad \frac{1}{n} \sum_{j=1}^n \widehat{\mathbf{d}}_j \mathbf{y}_j^\top \xrightarrow{p} \mathbb{E}[\mathbf{d}\mathbf{y}^\top] \quad (120)$$

Consequently, the empirical criterion

$$\widehat{\text{FC}} := \left\| \frac{1}{n} \sum_{j=1}^n \widehat{\mathbf{d}}_j \mathbf{y}_j^\top \right\|_F - \left\| \frac{1}{n} \sum_{j=1}^n \widehat{\mathbf{e}}_j \mathbf{x}_j^\top \right\|_F \quad (121)$$

is a consistent estimator of the population criterion:

$$\widehat{\text{FC}} \xrightarrow{p} \left\| \mathbb{E}[\mathbf{d}\mathbf{y}^\top] \right\|_F - \left\| \mathbb{E}[\mathbf{e}\mathbf{x}^\top] \right\|_F =: \text{FC} \quad (122)$$

E.3. Algorithms for estimating the two criteria

In this section, the observations $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times n}$ lie in two dimensions: the m views, and the n independent samples. We use the index i for the views, and k for the samples.

Algorithm 1 Likelihood-based criterion

Input: Observations $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times n}$ corresponding to two variables, residuals $\mathbf{E}, \mathbf{D} \in \mathbb{R}^{m \times n}$ corresponding to the directions $x \rightarrow y$ and $y \rightarrow x$, respectively.

1. Compute the covariance matrices for the direction $x \rightarrow y$:

$$\widehat{\Sigma}_x = \frac{1}{n} \sum_{j=1}^n \mathbf{X}_j \mathbf{X}_j^\top \quad \text{and} \quad \widehat{\Sigma}_e = \frac{1}{n} \sum_{j=1}^n \mathbf{E}_j \mathbf{E}_j^\top \quad (123)$$

and for the direction $y \rightarrow x$:

$$\widehat{\Sigma}_y = \frac{1}{n} \sum_{j=1}^n \mathbf{Y}_j \mathbf{Y}_j^\top \quad \text{and} \quad \widehat{\Sigma}_d = \frac{1}{n} \sum_{j=1}^n \mathbf{D}_j \mathbf{D}_j^\top . \quad (124)$$

2. Compute the likelihood based-criterion:

$$\widehat{\text{LR}} = -\log |\widehat{\Sigma}_x| - \log |\widehat{\Sigma}_e| + \log |\widehat{\Sigma}_y| + \log |\widehat{\Sigma}_d| . \quad (125)$$

Output: Likelihood-based criterion $\widehat{\text{LR}}$.

Algorithm 2 Cross-covariance-based criterion

Input: Observations $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{m \times n}$ corresponding to two variables, residuals $\mathbf{E}, \mathbf{D} \in \mathbb{R}^{m \times n}$ corresponding to the directions $x \rightarrow y$ and $y \rightarrow x$, respectively.

1. Compute the second-moment matrices

$$\mathbf{M}_1 = \frac{1}{n} \sum_{j=1}^n \mathbf{E}_j \mathbf{X}_j^\top \quad \text{and} \quad \mathbf{M}_2 = \frac{1}{n} \sum_{j=1}^n \mathbf{D}_j \mathbf{Y}_j^\top . \quad (126)$$

2. Compute the cross-covariance-based criterion:

$$\widehat{\text{FC}} = \|\mathbf{M}_2\|_F - \|\mathbf{M}_1\|_F . \quad (127)$$

Output: Cross-covariance-based criterion $\widehat{\text{FC}}$.

E.4. How to find the root variable

Assume we have a scalar criterion of the causal direction between two variables $\mathbf{x}_j = (x_j^1, \dots, x_j^m)^\top$ and $\mathbf{x}_k = (x_k^1, \dots, x_k^m)^\top$, such that the criterion is positive if $\mathbf{x}_j \rightarrow \mathbf{x}_k$ and negative in the opposite direction. Let M_{jk} denote the criterion obtained for the pair (j, k) , and collect all pairwise criteria in a matrix \mathbf{M} with a zero diagonal. We now have a matrix \mathbf{M} whose entries determine the causal direction by their sign: $M_{ij} > 0$ indicates \mathbf{x}_i causes \mathbf{x}_j , and $M_{ij} < 0$ indicates that \mathbf{x}_j causes \mathbf{x}_i . In this setting, \mathbf{x}_j is a root variable if and only if the j -th row of \mathbf{M} contains only non-negative entries. We follow Hyvärinen and Smith (2013, Section 3.2.2) to derive a principled way for finding the root variable given \mathbf{M} . We next recall their argument.

Suppose that the entries in \mathbf{M} follow a normal distribution, so that $M_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma^2)$ where σ^2 models the estimation error due to finite samples (and is supposed to be constant for simplicity). Then, the log-likelihood of \mathbf{x}_i being the root

variable is

$$\log \prod_{j=1}^p \mathbb{P}(\mu_{ij} > 0 | M_{ij}) = \sum_{j=1}^p \log \mathbb{P}\left(\frac{\mu_{ij} - M_{ij}}{\sigma} > \frac{-M_{ij}}{\sigma} \mid M_{ij}\right) \quad (128)$$

$$= \sum_{j=1}^p \log \Phi\left(\frac{M_{ij}}{\sigma}\right) \approx -\frac{1}{2\sigma^2} \sum_{j=1}^p \min(0, M_{ij})^2 \quad (129)$$

where Φ is the cumulative distribution function of a standardized Gaussian, whose log we then approximated. We therefore obtain the index of the root variable using

$$\arg \min_i \sum_{j=1}^p \min(0, M_{ij})^2. \quad (130)$$

E.5. One-step Feasible Generalized Least Squares (FGLS)

The FGLS procedure begins with an Ordinary Least Squares (OLS) regression in each view to obtain residuals (e_j^1, \dots, e_j^m) , which are then used to estimate the cross-view covariance of the disturbances. In a second step, the OLS coefficients are adjusted via Generalized Least Squares using this estimated covariance. In the population limit, this estimator coincides with OLS but is asymptotically more efficient whenever disturbances are correlated across views, as is the case here.

The estimation proceeds as follows. Assume we have observed n independent samples. Here, we consider the model

$$\mathbf{x}^i = \mathbf{T}^i \mathbf{x}^i + \mathbf{e}^i, \quad i \in \llbracket 1, m \rrbracket \quad (131)$$

where \mathbf{x}^i and \mathbf{e}^i are the reordered observations and disturbances, respectively, and \mathbf{T}^i are strictly lower triangular matrices. By construction, each variable x_j^i only depends on its predecessors $\mathbf{x}_{<j}^i := (x_1^i, \dots, x_{j-1}^i)$.

Fix a variable index $j \in \llbracket 1, p \rrbracket$. In view i , the structural equation reads

$$x_j^i = \sum_{k=1}^{j-1} T_{jk}^i x_k^i + e_j^i. \quad (132)$$

Stacking all views, we obtain the block regression model

$$\mathbf{x}'_j = \mathbf{X}_{<j} \mathbf{T}_{<j} + \mathbf{e}'_j \quad (133)$$

where

- $\mathbf{x}'_j = (x_j^1, \dots, x_j^m)^\top \in \mathbb{R}^{mn}$ stacks all samples across the m views,
- $\mathbf{X}_{<j} = \text{diag}(\mathbf{x}_{<j}^1, \dots, \mathbf{x}_{<j}^m) \in \mathbb{R}^{mn \times m(j-1)}$ is block-diagonal and contains the predecessors,
- $\mathbf{T}_{<j} = (\mathbf{T}_{j,<j}^1, \dots, \mathbf{T}_{j,<j}^m)^\top \in \mathbb{R}^{m(j-1)}$ collects the coefficients,
- $\mathbf{e}'_j = (e_j^1, \dots, e_j^m)^\top \in \mathbb{R}^{mn}$ are disturbances.

The disturbance vector \mathbf{e}'_j has covariance

$$\Omega_j := \mathbb{E}[\mathbf{e}'_j (\mathbf{e}'_j)^\top] = \Sigma_{\mathbf{e}_j} \otimes \mathbf{I}_n \in \mathbb{R}^{mn \times mn} \quad (134)$$

where $\Sigma_{\mathbf{e}_j} \in \mathbb{R}^{m \times m}$ captures cross-view covariances for component j .

The one-step FGLS procedure is:

1. *OLS step:*

$$\widehat{\mathbf{T}}_{<j}^{\text{OLS}} = (\mathbf{X}_{<j}^\top \mathbf{X}_{<j})^{-1} \mathbf{X}_{<j}^\top \mathbf{x}'_j \quad (135)$$

yielding residuals $\hat{e}'_j = \mathbf{x}'_j - \mathbf{X}_{<j} \widehat{\mathbf{T}}_{<j}^{\text{OLS}}$.

2. *Covariance estimation:*

$$\widehat{\Sigma}_{e_j} = \frac{1}{n} \sum_{k=1}^n \hat{e}'_{j,k} (\hat{e}'_{j,k})^\top \quad (136)$$

where k denotes the sample index.

3. *FGLS update:*

$$\widehat{\mathbf{T}}_{<j}^{\text{FGLS}} = (\mathbf{X}_{<j}^\top \widehat{\Sigma}_{e_j}^{-1} \mathbf{X}_{<j})^{-1} \mathbf{X}_{<j}^\top \widehat{\Sigma}_{e_j}^{-1} \mathbf{x}'_j . \quad (137)$$

This estimator reduces to OLS if Σ_{e_j} is diagonal (no cross-view correlation). More generally, it achieves asymptotic efficiency by exploiting cross-view correlations.

E.6. Full algorithm

Algorithm 3 Pairwise-comparison algorithms for multi-view causal discovery

Input: Observations $\mathbf{X} \in \mathbb{R}^{m \times p \times n}$.

1. *Preprocess the data \mathbf{X} .* Standardize observations over the samples axis.

2. *Estimate the causal ordering \mathbf{P} .*

- (a) Perform pairwise regressions between each pair of variables, \mathbf{x}_j on \mathbf{x}_i , such that $i \neq j$. Compute the residuals $\mathbf{e}_{i \rightarrow j}$.
- (b) Compute a skew-symmetric matrix $\mathbf{M} \in \mathbb{R}^{p \times p}$ that determines the causal direction between pairs of variables.

The entries M_{ij} are computed with Eq. 5 or Eq. 6 for all $i < j$, using the regression results $(\mathbf{x}_i, \mathbf{x}_j, \mathbf{e}_{i \rightarrow j}, \mathbf{e}_{j \rightarrow i})$. Then set $M_{ji} = -M_{ij}$. M_{ij} positive means \mathbf{x}_i causes \mathbf{x}_j , negative means \mathbf{x}_j causes \mathbf{x}_i .

- (c) Determine the root variable k : it causes all others, so M_{kj} should be positive for all j .

$$k = \underset{i}{\operatorname{argmin}} \sum_j \min(0, M_{ij})^2 . \quad (138)$$

- (d) Remove the root variable k and replace observations with residuals obtained when regressing on k . Repeat.
- (e) Store the estimated ordering in a permutation matrix \mathbf{P} , and reorder the variables in \mathbf{X} according to this permutation. This yields a new model $\mathbf{x}'^i = \mathbf{T}^i \mathbf{x}'^i + \mathbf{e}'^i$, where \mathbf{x}'^i and \mathbf{e}'^i are reordered versions of \mathbf{x}^i and \mathbf{e}^i , respectively, and the \mathbf{T}^i are strictly lower triangular matrices.

3. *Estimate the adjacency matrices \mathbf{B}^i .* For each variable $j \in \llbracket 2, p \rrbracket$, estimate the j -th row of all matrices \mathbf{T}^i with one-step Feasible Generalized Least Squares (Zellner, 1962). Recover matrices $\mathbf{B}^i = \mathbf{P}^\top \mathbf{T}^i \mathbf{P}$.

Output: Adjacency matrices $(\mathbf{B}^1, \dots, \mathbf{B}^m) \in \mathbb{R}^{m \times p \times p}$.

E.7. Computational complexity

The computational complexity of Algorithm 3 can be broken down across the different steps.

In Step 1, the preprocessing is in $O(m \cdot p \cdot n)$.

In Step 2.a., we solve $O(p^2)$ LS regression problems, each of which has complexity $O(m \cdot n)$.

1650 In Step 2.b., we compute $O(p^2)$ entries of a matrix. Each entry is computed either using the cross-covariance-based
 1651 criterion in Algorithm 2 which is in $O(m^2 \cdot n)$, or using the likelihood-based criterion in Algorithm 1 which is in $O(m^2 \cdot$
 1652 $n + m^3)$ due to computing log determinants. So the total complexity of this step is $O(m^2 \cdot p^2 \cdot n)$ or $O(m^2 \cdot p^2 \cdot n + m^3 \cdot p^2)$.

1653 In Step 2.c., we make $O(p^2)$ operations.
 1654

1655 In Step 2.d., we simplify/redefine a matrix.
 1656

1657 These steps 2.a-d., are repeated $O(p)$ times, until all variables are removed, so this yields a complexity in either $O(m^2 \cdot p^3 \cdot n)$
 1658 or $O(m^2 \cdot p^3 \cdot n + m^3 \cdot p^3)$.

1659 In Step 2.e., we reorder a matrix, which costs $O(m \cdot p \cdot n)$ operations.
 1660

1661 In Step 3, we run the one-step Feasible GLS procedure, which is in $O(m^3 \cdot p^3 \cdot n)$. It scales cubically in the m views and
 1662 p dimensions as it involves matrix multiplications and inversions. It then does m matrix multiplications. Each of them
 1663 involves permutation matrices which has a quadratic cost $O(p^2)$. So the cost of this step is dominated by $O(m^3 \cdot p^3 \cdot n)$.

1664 The total computational complexity of the algorithm is thus $O(m^3 \cdot p^3 \cdot n)$.
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673
 1674
 1675
 1676
 1677
 1678
 1679
 1680
 1681
 1682
 1683
 1684
 1685
 1686
 1687
 1688
 1689
 1690
 1691
 1692
 1693
 1694
 1695
 1696
 1697
 1698
 1699
 1700
 1701
 1702
 1703
 1704

F. Identifiability of LiMVAM with shared disturbances

In this section, we analyze the identifiability of the LiMVAM model with shared disturbances

$$\mathbf{x}^i = \mathbf{B}^i \mathbf{x}^i + \mathbf{D}^i \mathbf{s} + \mathbf{n}^i, \quad i \in \llbracket 1, m \rrbracket \quad (139)$$

where the \mathbf{D}^i are diagonal matrices with positive entries on the diagonal, \mathbf{s} has mutually independent entries and second-order moment $\mathbb{E}[\mathbf{s}\mathbf{s}^\top] = \mathbf{I}_p$, the view-specific noises are Gaussian $\mathbf{n}^i \sim \mathcal{N}(\mathbf{0}, \Sigma^i)$ with diagonal Σ^i , and the vectors \mathbf{s} and $\mathbf{n}^i, i = 1, \dots, m$, are mutually independent.

F.1. Proof of Theorem 4 (first claim)

Consider two sets $\Theta = (\mathbf{D}^1, \dots, \mathbf{D}^m, \Sigma^1, \dots, \Sigma^m, \mathbf{B}^1, \dots, \mathbf{B}^m)$ and $\Theta' = (\mathbf{D}'^1, \dots, \mathbf{D}'^m, \Sigma'^1, \dots, \Sigma'^m, \mathbf{B}'^1, \dots, \mathbf{B}'^m)$ that parameterize the same statistical model in Eq. 139, with \mathbf{D}^i and \mathbf{D}'^i being diagonal matrices with positive entries on the diagonal, Σ^i and Σ'^i being diagonal matrices with non-zero entries on the diagonal, and \mathbf{B}^i and \mathbf{B}'^i being DAG matrices. The model in Eq. 139 can be reformulated as

$$\mathbf{x}^i = (\mathbf{I} - \mathbf{B}^i)^{-1} \mathbf{D}^i (\mathbf{s} + (\mathbf{D}^i)^{-1} \mathbf{n}^i) \quad (140)$$

which corresponds to the Shared ICA model

$$\mathbf{x}^i = \mathbf{A}^i (\mathbf{s} + \tilde{\mathbf{n}}^i) \quad (141)$$

for $\mathbf{A}^i = (\mathbf{I} - \mathbf{B}^i)^{-1} \mathbf{D}^i$ and $\tilde{\mathbf{n}}^i = (\mathbf{D}^i)^{-1} \mathbf{n}^i$. We can observe that the unmixing matrices $\mathbf{W}^i = (\mathbf{A}^i)^{-1} = (\mathbf{D}^i)^{-1} (\mathbf{I} - \mathbf{B}^i)$ and $\mathbf{W}'^i = (\mathbf{A}'^i)^{-1} = (\mathbf{D}'^i)^{-1} (\mathbf{I} - \mathbf{B}'^i)$ belong to the domain \mathcal{W} , which is included in the space of invertible matrices. Thus, the two sets of \mathbf{W}^i and \mathbf{W}'^i matrices are valid sets of unmixing matrices for the same multi-view ICA model in Eq. 141. Moreover, Assumption 5 allows the rescaled noises $\tilde{\mathbf{n}}^i$ to meet the noise diversity condition of Shared ICA. So, from the identifiability theory of multi-view ICA (Richard et al., 2021, Theorem 1), we know that there exist a sign-permutation matrix \mathbf{Q} such that for any view $i \in \llbracket 1, m \rrbracket$, we have

$$\begin{aligned} \mathbf{W}'^i &= \mathbf{Q}^\top \mathbf{W}^i \\ \Sigma'^i &= \mathbf{Q}^\top \Sigma^i \mathbf{Q} . \end{aligned} \quad (142)$$

Note that, contrary to the single-view context, the fact that we obtain \mathbf{W} up to *sign* and permutation rather than *scale* and permutation is because $\mathbb{E}[\mathbf{s}\mathbf{s}^\top] = \mathbf{I}$. Then, we apply Lemma 5, which shows that being in the domain \mathcal{W} imposes $\mathbf{Q} = \mathbf{I}$, $\mathbf{D}'^i = \mathbf{D}^i$, and $\mathbf{B}'^i = \mathbf{B}^i$.

F.2. Proof of Theorem 4 (second claim)

Here we prove that, under Assumptions 5 and 3, the decomposition of matrices \mathbf{B}^i into matrices \mathbf{T}^i and \mathbf{P} is unique.

In particular, Assumption 3 requires the shared ordering to be total, *i.e.* there exists a directed path between any pair of variables of the (union) graph. This assumption can be reformulated mathematically, as follows. Let $\bar{\mathbf{P}}$ be the permutation matrix of the shared ordering ($\bar{\mathbf{P}}$ is unique since the ordering is total) and $\bar{\mathbf{T}}^i = \bar{\mathbf{P}} \mathbf{B}^i \bar{\mathbf{P}}^\top$ be the corresponding strictly lower triangular matrices. We define

$$\bar{\mathbf{T}}^U = \sum_{i=1}^m \text{abs}(\bar{\mathbf{T}}^i) \quad (143)$$

which is strictly lower triangular and contains the causal weights of all views, where abs denotes the element-wise absolute value function. We know from graph theory that, for the adjacency matrix \mathbf{B} of a directed graph and any integer $k \geq 1$, the matrix \mathbf{B}^k counts walks of length k , in the sense that $\mathbf{B}_{ij}^k \neq 0$ iff there is a path of length k that goes from variable j to variable i . So, the fact that the shared ordering is total implies that, for any pair of variables $i > j$, there exists an integer k such that $(\bar{\mathbf{T}}^U)_{ij}^k \neq 0$. Equivalently, the strictly lower triangular part of $\sum_{k=1}^{p-1} (\bar{\mathbf{T}}^U)^k$ only has non-zero elements.

Now, consider the two sets $\Theta = (\Sigma^1, \dots, \Sigma^m, \mathbf{P}, \mathbf{T}^1, \dots, \mathbf{T}^m)$ and $\Theta' = (\bar{\Sigma}^1, \dots, \bar{\Sigma}^m, \bar{\mathbf{P}}, \bar{\mathbf{T}}^1, \dots, \bar{\mathbf{T}}^m)$ that parameterize the same statistical model in Eq. 139, where \mathbf{P} is some permutation matrix, \mathbf{T}^i are strictly lower triangular matrices, and $\bar{\mathbf{P}}$ and $\bar{\mathbf{T}}^i$ are the particular matrices given by Assumption 3. Note that no assumption on the sparsity of the \mathbf{T}^i is made. From Theorem 4, we know that $\Sigma^i = \bar{\Sigma}^i$.

1760 Let us define

$$1761 \bar{\mathbf{W}}^U = \sum_{i=1}^m \text{abs}(\bar{\mathbf{W}}^i) \quad (144)$$

1765 where $\bar{\mathbf{W}}^i = \mathbf{I} - \bar{\mathbf{P}}^\top \bar{\mathbf{T}}^i \bar{\mathbf{P}}$. The non-zero elements of $\bar{\mathbf{P}}^\top \bar{\mathbf{T}}^i \bar{\mathbf{P}}$ are outside of the diagonal, so matrices \mathbf{I} and $\bar{\mathbf{P}}^\top \bar{\mathbf{T}}^i \bar{\mathbf{P}}$ contain non-zero elements at different locations. Thus, we have

$$1766 \text{abs}(\mathbf{I} - \bar{\mathbf{P}}^\top \bar{\mathbf{T}}^i \bar{\mathbf{P}}) = \mathbf{I} + \text{abs}(\bar{\mathbf{P}}^\top \bar{\mathbf{T}}^i \bar{\mathbf{P}}) . \quad (145)$$

1769 Furthermore, applying $\bar{\mathbf{P}}$ to the rows and columns of $\bar{\mathbf{T}}^i$ only shuffles its entries, without modifying their values. So we have

$$1772 \text{abs}(\bar{\mathbf{P}}^\top \bar{\mathbf{T}}^i \bar{\mathbf{P}}) = \bar{\mathbf{P}}^\top \text{abs}(\bar{\mathbf{T}}^i) \bar{\mathbf{P}} . \quad (146)$$

1774 Consequently,

$$1775 \bar{\mathbf{W}}^U = \sum_{i=1}^m \mathbf{I} + \bar{\mathbf{P}}^\top \text{abs}(\bar{\mathbf{T}}^i) \bar{\mathbf{P}} = m\mathbf{I} + \bar{\mathbf{P}}^\top \bar{\mathbf{T}}^U \bar{\mathbf{P}} . \quad (147)$$

1779 Next, we apply the same reasoning to the alternative set of parameters, given by $\mathbf{W}^i = \mathbf{I} - \mathbf{P}^\top \mathbf{T}^i \mathbf{P}$, and we consider $\mathbf{W}^U = \sum_{i=1}^m \text{abs}(\mathbf{W}^i)$. The proof of Theorem 4 already implied that $\bar{\mathbf{W}}^i = \mathbf{W}^i$ for all i . Thus, we have $\bar{\mathbf{W}}^U = \mathbf{W}^U$ and

$$1782 \bar{\mathbf{P}}^\top \bar{\mathbf{T}}^U \bar{\mathbf{P}} = \mathbf{P}^\top \mathbf{T}^U \mathbf{P} \quad (148)$$

1783 where \mathbf{T}^U is defined in a similar way as $\bar{\mathbf{T}}^U$, except that its strictly lower triangular part can be sparse. Raising this equation to power k , for $k = 1, \dots, p-1$, we get

$$1784 \bar{\mathbf{P}}^\top \left(\sum_{k=1}^{p-1} (\bar{\mathbf{T}}^U)^k \right) \bar{\mathbf{P}} = \mathbf{P}^\top \left(\sum_{k=1}^{p-1} (\mathbf{T}^U)^k \right) \mathbf{P} \quad (149)$$

1790 where we recall that the strictly lower triangular part of $\sum_{k=1}^{p-1} (\bar{\mathbf{T}}^U)^k$ only has non-zero elements. Using Lemma 6 on Eq. 149, we obtain that $\mathbf{P} = \bar{\mathbf{P}}$, and thus $\mathbf{T}^i = \bar{\mathbf{T}}^i$ for all i . In conclusion, all the sets of DAG decompositions that parameterize the same model are equal. We conclude that matrices \mathbf{P} and \mathbf{T}^i are unique, and thus identifiable in our terminology.

1795 F.3. Results for view-specific causal orderings

1796 Next, we consider a case which is outside of the theory of the main paper, although a simple extension: we allow the causal orderings to be different in different views. It turns out that in the view-specific \mathbf{P}^i case, identifiability is obtained by assuming that the directed acyclic graph \mathbf{B}^i is dense enough in each view, as formalized in the following assumption. However, since the \mathbf{B}^i can be permuted to strictly lower triangular matrices, they cannot be denser than having $\frac{p(p-1)}{2}$ non-zero entries.

1802 **Assumption 7** (Dense connectivity in each view) *For each view i , the matrix \mathbf{B}^i has exactly $\frac{p(p-1)}{2}$ non-zero entries.*

1804 Using Assumption 7, the following theorem states that, in addition to identifying \mathbf{B}^i , \mathbf{D}^i , and Σ^i , one can also identify \mathbf{T}^i and \mathbf{P}^i .

1806 **Theorem 10** (Identifiability of multiple causal orderings) *Consider the LiMVAM model with multiple causal orderings. Consider the quantities $(\mathbf{D}^1, \dots, \mathbf{D}^m, \Sigma^1, \dots, \Sigma^m, \mathbf{P}^1, \dots, \mathbf{P}^m, \mathbf{T}^1, \dots, \mathbf{T}^m)$ as the set of parameters to be estimated. Under Assumptions 5 and 7, all these parameters are identifiable.*

1810 In the context of view-specific causal orderings \mathbf{P}^i , we prove that, under Assumptions 5 and 7, the decomposition of matrices \mathbf{B}^i into matrices \mathbf{T}^i and \mathbf{P}^i is unique.

1812 Consider two sets of parameters $\Theta = (\Sigma^1, \dots, \Sigma^m, \mathbf{P}^1, \dots, \mathbf{P}^m, \mathbf{T}^1, \dots, \mathbf{T}^m)$ and $\Theta' = (\Sigma'^1, \dots, \Sigma'^m, \mathbf{P}'^1, \dots, \mathbf{P}'^m, \mathbf{T}'^1, \dots, \mathbf{T}'^m)$ that parameterize the same statistical model in Eq. 139, where $\mathbf{P}^i, \mathbf{P}'^i$ are permutation

matrices, and T^i, T'^i are strictly lower triangular matrices. Note that here we parameterize by P^i and T^i rather than B^i or W^i . From Theorem 4, we know that $\Sigma^i = \Sigma'^i$ and that the resulting causal matrices must be equal:

$$(P^i)^\top T^i P^i = (P'^i)^\top T'^i P'^i . \quad (150)$$

Assumption 7 states that, for each view i , the matrix $B^i = (P^i)^\top T^i P^i = (P'^i)^\top T'^i P'^i$ contains exactly $\frac{p(p-1)}{2}$ non-zero elements, so it is also the case for T^i and T'^i which thus represent fully connected graphs. So, from Lemma 6, we deduce that, for each view i , we have $P^i = P'^i$ and $T^i = T'^i$, which concludes the proof.

G. ICA-based Algorithm for LiMVAM with shared disturbances

G.1. ICA-based Algorithm

Algorithm 4 ICA-LiMVAM

1. *Estimate the adjacency matrices B^i .*

- (a) Estimate the unmixing matrices W^i and the noise variance matrices Σ^i , by running the Shared ICA estimation algorithm on the data.

The algorithm returns an estimate of the true unmixing matrix $W^i = M(I - B^i)$ but up to sign-permutation matrix M that is the same for all views (Richard et al., 2021).

- (b) Determine the sign-permutation indeterminacy M , using the structure of the underlying $(A^i)^{-1}$ that has a diagonal of ones. Thus, we can adapt the simple two-step heuristic by Shimizu et al. (2006).

First, find a permutation matrix M such that MW has a non-zero diagonal, by solving with the Hungarian algorithm

$$\min_M \sum_{j=1}^p \frac{1}{|(MW)_{jj}|}, \quad W = \sum_{i=1}^m \text{abs}(W^i), \quad (151)$$

and then rescale the rows of M to ensure that MW has a diagonal of ones

$$M_{ij} \leftarrow \text{sign}((MW)_{ii})M_{ij}, \quad (152)$$

where the sign function outputs 1 or -1 .

- (c) Determine the scale matrices $D^i = \text{diag}\left(\frac{1}{(MW^i)_{11}}, \dots, \frac{1}{(MW^i)_{pp}}\right)$.

- (d) Determine the causal matrices $B^i = I - D^i M W^i$ and update the noise variance matrices with $\Sigma^i \leftarrow (D^i)^2 \Sigma^i$.

2. *Determine the causal ordering P .*

The DAG decomposition states that $B^i := P^\top T^i P$ where T^i is lower triangular, so we penalize for that

$$\min_P \sum_{l \geq k} (P B P^\top)_{kl}^2, \quad B = \sum_{i=1}^m \text{abs}(B^i). \quad (153)$$

We approximately minimize this objective using the heuristic algorithm presented in Shimizu et al. (2006, Algorithm C).

Note that the algorithm described in Chen et al. (2024) finds the same adjacency matrices as we do. This is because of the sign-permutation matrix M from ICA: to determine that matrix, we ensure MW has a diagonal of ones by dividing by the scalings D^i . For (Chen et al., 2024), these scalings are not part of the model: they are simply a part of the algorithm which determines M . For us, these scalings are part of our model.

G.2. Computational complexity

We next detail the worst-case computational complexity of the ICA-LiMVAM algorithm which is in $O(tnmp^3 + mp^5)$, for t iterations of the SharedICA-ML algorithm, n samples, m views and p components. Note that the term in p^5 comes from Algorithm C of the original LiNGAM and is worst-case: in practice, it took a couple of iterations only in our experiments. Below are the detail of the derivations for each step of Algorithm 4.

Computational complexity The computational complexity of ICA-LiMVAM can be broken down across the four steps in Algorithm 4. Step 1 calls SharedICA, with its most costly variant (SharedICA-ML) running in $O(t \cdot n \cdot m \cdot p^3)$, where t is the number of iterations, n the number of samples, m the number of views, and p the number of components. Step 2

solves the permutation indeterminacy via a linear sum assignment in $O(m \cdot p^3)$, instead of trying $p!$ permutations. Step 3 involves m multiplications with permutation matrices and additions, costing $O(mp^2)$. Step 4 optimizes over permutations. Instead of trying $p!$ permutations, we use the heuristic Algorithm C from LiNGAM. It is $O(m \cdot p^5)$ in the worst case (and much lower in practice). Further details are provided next:

1. Run the ICA algorithm SharedICA-ML, $O(tnmp^3)$

Each of the t iterations of the optimization performs the so-called “E-step” and “M-step” of the E-M algorithm, detailed in Section 4 of Richard et al. (2021). Their complexity is in mnp^3 : the intuition is that for each view m and each sample n require some matrix operations (*e.g.* addition, multiplication, computing gradients and Hessians) all of which are dominated by p^3 . So the final complexity is in $O(tnmp^3)$.

2. Solve the permutation-sign indeterminacy, $O(mp^3)$

Computing the objective can be done in $O(p^2)$. The detail of this is: we begin by computing the matrix W which is in $O(mp^2)$, then multiplying it with a permutation matrix which is in $O(p^2)$, and finally summing the diagonal elements which is in $O(p)$.

Finding the correct permutation matrix can be done by solving a linear-sum assignment problem; this can be done using the Hungarian algorithm which is in $O(p^3)$.

3. Update the entries of a matrix which is in $O(p^2)$.

4. Update the entries of a matrix which is in $O(p^2)$.

Adding the complexities thus far leads to a complexity of $O(mp^2 + p^3)$, which is dominated by $O(mp^3)$, for Step 2.

5. Compute the causal matrices, $O(mp^2)$

Each of the m causal matrices requires multiplying a dense matrix with a permutation which is in $O(p^2)$ and then adding a matrix which is in $O(p^2)$. The final cost is $O(mp^2)$.

6. Find the causal ordering(s), $O(mp^5)$

Algorithm B in the original LiNGAM paper finds a row that has all zeros $O(p^2)$, removes it, does this p times. Its complexity is in $O(p^3)$.

Algorithm C in the original LiNGAM paper calls Algorithm B at most $p(p - 1)/2 = O(p^2)$ times, so its final complexity is $O(p^2 p^3) = O(p^5)$.

When the causal ordering is shared across views: first we compute sum m matrices with p^2 entries which is in $O(mp^2)$, and then we use Algorithm C once which is in $O(p^5)$. The final complexity is in $O(mp^2 + p^5)$ which is dominated by $O(mp^5)$.

When the causal ordering is view-dependent: for each of the m views, we use Algorithm C which is in $O(p^5)$, so the final complexity is in $O(mp^5)$.

H. Experiments

H.1. Synthetic experiments

H.1.1. DATA GENERATION IN FIGURE 2

The causal effect matrices $B^i = P^\top T^i P$ are generated from a random permutation P and strictly lower triangular T^i obtained from a standard Gaussian; the diagonal matrices D^i are drawn from a uniform density between 0.1 and 3; the common disturbances in s can be either Gaussian or non-Gaussian: Gaussian disturbances s_j are generated from a standardized Gaussian and their corresponding noises n_j^i have standard deviation Σ_{jj}^i , obtained by sampling from a uniform density between 0 and 1, while non-Gaussian disturbances are generated from a Laplace distribution (with scale parameter equal to $\frac{1}{2}$) and their corresponding noises all have a std of $\frac{1}{2}$. We use $m = 5$ views, $p = 4$ variables, and vary the number of samples n between 10^2 and 10^4 .

H.1.2. SIMULATION STUDY IN HIGHER DIMENSION

We evaluate the methods by measuring the estimation error on the matrices B^i across a range of settings involving more views and components than those considered in the main text. Specifically, we vary the number of views among $\{3, 5, 8, 12, 16, 20\}$ and the number of components among $\{3, 6, 9, 12\}$, using 1000 samples and 50 random seeds for each configuration.

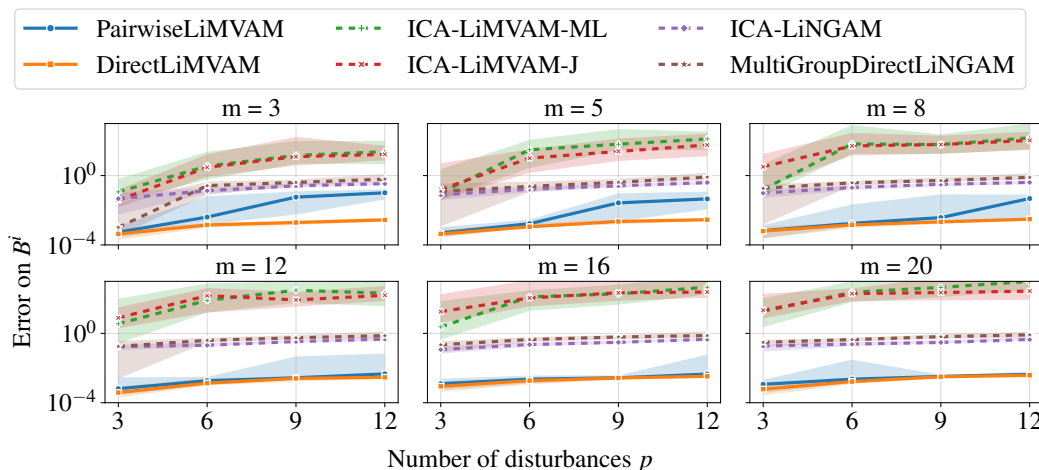


Figure 6. Separation performance of three pairwise-comparisons-based algorithms — MultiGroupDirectLiNGAM from Shimizu (2012), and DirectLiMVAM and PairwiseLiMVAM which use the criteria in Eq. 5 and Eq. 6, respectively — and three ICA-based algorithms — ICA-LiNGAM, which naively uses Shimizu et al. (2006) separately in each view, and our implementation of the two variants of ICA-LiMVAM (Chen et al., 2024). We varied the number of views and components. Disturbances are generated using the generalized normal distribution with shape parameter $\beta \in \{1.5, 2, 2.5\}$. The metric used is the ℓ_2 -distance between true and estimated causal effect matrices B^i (lower is better).

Data are generated from the model in Eq. 3, where disturbances e_j are split into three types—sub-Gaussian, Gaussian, and super-Gaussian—which motivates the choice of disturbance counts as multiples of 3. Sub-Gaussian disturbances follow a generalized normal distribution with shape parameter $\beta = 2.5$; super-Gaussian disturbances use $\beta = 1.5$; Gaussian disturbances correspond to $\beta = 2$. Each T^i is sampled as a strictly lower-triangular matrix with Gaussian entries, and a common permutation matrix P is used to construct the $B^i = P^\top T^i P$. For each method, we report the median error on the matrices B^i , with the 25th and 75th percentiles shown as error bars.

As shown in Fig. 6, DirectLiMVAM and PairwiseLiMVAM consistently achieve lower errors than the other approaches. DirectLiMVAM performs best when the number of views is limited (3, 5, or 8) while the number of components remains relatively high (9 or 12), thereby outperforming PairwiseLiMVAM in these settings. This slightly surprising observation that DirectLiNGAM is a bit better than PairwiseLiNGAM is perhaps due to the fact that there the particular non-Gaussianities used here are ill-suited for the Gaussian likelihood underlying PairwiseLiNGAM.

In contrast, ICA-LiNGAM and MultiGroupDirectLiNGAM achieve slightly better-than-average performance when both the number of views and the number of components are small, but their performance quickly degrades to the average

level as dimensionality increases. This behavior reflects their inability to properly handle Gaussian disturbances. Finally, both ICA-LiMVAM variants fail completely, as they are not designed to operate without shared disturbances. Overall, the estimation error tends to increase with the number of components.

H.1.3. COMPARISON WITH A MULTI-DOMAIN METHOD

Here, we illustrate the advantage of using multi-view methods like ours when cross-view correlations are present in the data. Specifically, we compare our DirectLiMVAM to the method of Perry et al. (2022), MSS (see Appendix I). We used their implementation with the KCI estimator, as this variant performs best in their paper. Since MSS recovers only the common DAG structure (and not the causal weights), we measure performance in terms of recovering the causal ordering.

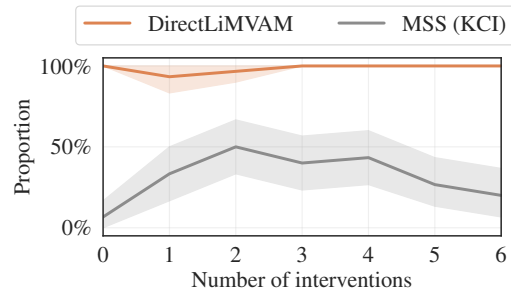


Figure 7. Proportion of runs in which the ordering is perfectly recovered (higher is better). The number of interventions corresponds to the number of entries in the \mathbf{B}^i that are allowed to vary across views. We compare our multi-view algorithm DirectLiMVAM to the multi-domain algorithm MSS of Perry et al. (2022) (using the KCI estimator).

We used the following experimental setup. Data were generated according to Eq. 1 with 6 views, 4 Gaussian disturbances, 500 samples, and the experiment was repeated over 30 random seeds. The noise variances were all fixed to one, and some entries in the matrices \mathbf{B}^i were allowed to vary across views. The number of such varying entries corresponds to the number of “interventions”, *i.e.* changes in causal mechanisms in the sense of Perry et al. (2022).

Figure 7 shows that DirectLiMVAM always recovers the correct ordering, regardless of how many entries in \mathbf{B}^i are allowed to vary. By contrast, MSS fails to recover the DAG when the number of interventions is either too small or too large (in line with their Figures 4 and 5), and even in the regime of sparse interventions it perfectly recovers the ordering in only about 50% of the runs. Finally, their method is substantially slower than ours (not shown in the graph).

These results illustrate the distinction between multi-view and multi-domain methods: by explicitly exploiting cross-view correlations, DirectLiMVAM can achieve stable and accurate recovery in settings

where a multi-domain method like MSS fails. Furthermore, this experiment confirms that our methods do not require changes in the \mathbf{B}^i as long as there are sufficient correlations between views, thereby extending the identifiability theories in the work just discussed.

H.1.4. TESTING ASSUMPTION 5

Assumption 5 provides a sufficient condition for the identifiability of the causal matrices \mathbf{B}^i in the LiMVAM model with shared, Gaussian disturbances. To evaluate the practical impact of this assumption, we simulated data from the model in Eq. 3, with shared disturbances as in Eq. 9. We used 5 views, 4 common disturbances in \mathfrak{s} (2 Gaussian and 2 Laplacian), and 1000 samples.

The scale matrices \mathbf{D}^i were drawn uniformly from the interval $[0.5, 2]$, and the strictly lower triangular matrices \mathbf{T}^i were sampled from a Gaussian distribution and then permuted to form the causal matrices \mathbf{B}^i . The noise variances Σ_{jj}^i corresponding to the Laplacian disturbances were fixed to $\frac{1}{2}$, while the noise variances corresponding to the two Gaussian disturbances, indexed by j and j' , were initially sampled uniformly in $[0, 1]$. To test the assumption, we selected an increasing number of views i such that the scaled noise variances $\frac{\Sigma_{jj}^i}{\mathbf{D}_{jj}^i} = \frac{\Sigma_{j'j'}^i}{\mathbf{D}_{j'j'}^i}$, ranging from 0 to 5. Note that Assumption 5 only fails when this condition holds across all 5 views. We repeated this experiment over 50 random seeds.

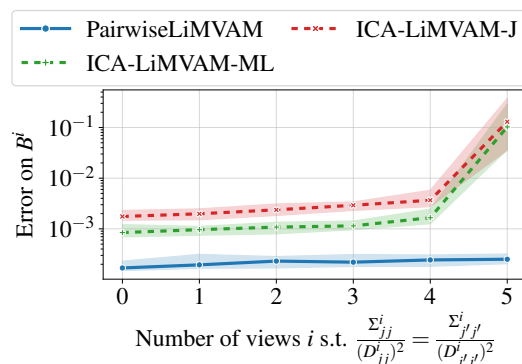


Figure 8. Effect of violating the noise diversity assumption (Assumption 5) on the estimation error of the causal matrices \mathbf{B}^i . We report the average ℓ_2 error over 50 repetitions for ICA-LiMVAM and PairwiseLiMVAM, as the number of views with identical scaled noise variances increases.

Figure 8 shows that ICA-LiMVAM maintains low estimation error on the \mathbf{B}^i matrices as long as the noise diversity

assumption holds, but its performance deteriorates sharply once the assumption is violated. This confirms the practical relevance of the assumption. Interestingly, PairwiseLiMVAM appears unaffected by the violation of this condition.

H.1.5. EFFECT OF THE ADJACENCY MATRICES' SPARSITY

Assumption 3 guarantees that the *shared causal ordering* \mathbf{P} and the \mathbf{T}^i can be uniquely recovered, as soon as there exists a directed path between any two variables in the graph. Mathematically, this corresponds to assuming that $\sum_{k=1}^{p-1} (\sum_{i=1}^m \text{abs}(\mathbf{T}^i))^k$ only has non-zero entries below the diagonal. In other words, it corresponds to assuming that the graph of $\sum_{i=1}^m \text{abs}(\mathbf{T}^i)$ is sufficiently dense and has a structure that leaves no indeterminacy between two variables.

In the following experiment, we consider a simpler version of this assumption and examine the performance of PairwiseLiMVAM and ICA-LiMVAM when varying the sparsity of the \mathbf{T}^i . Intuitively, if the \mathbf{T}^i are dense enough, then Assumption 3 will automatically be fulfilled.

In the experiment, we also study the case of *view-specific (or multiple) causal orderings* — introduced in Appendix F.3 — which is outside of the theory of the paper. When causal orderings are not shared across views, assuming that all \mathbf{B}^i are dense (in the sense that each \mathbf{B}^i has exactly $\frac{p(p-1)}{2}$ non-zero entries) is sufficient to make the (multiple) orderings \mathbf{P}^i and the \mathbf{T}^i identifiable. In Appendix F.3, this assumption is named Assumption 7.

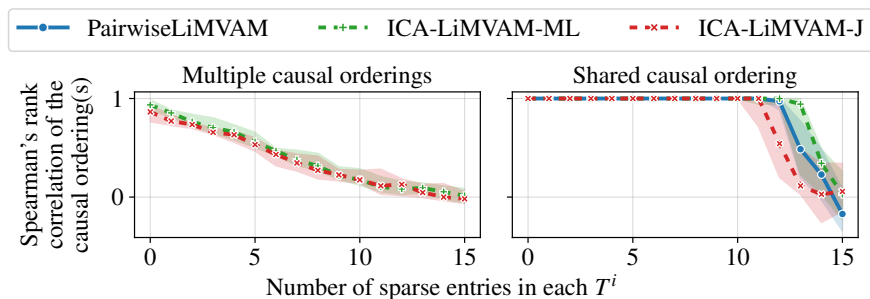


Figure 9. Spearman’s rank correlation between true and estimated causal orderings as a function of the number of sparse entries in each \mathbf{T}^i . Left: In the multiple-ordering setting, performance of ICA-LiMVAM-ML and ICA-LiMVAM-J degrades linearly with increasing sparsity, as the assumption of full support is violated. PairwiseLiMVAM is excluded since it assumes a shared ordering. Right: In the shared-ordering setting, all methods maintain near-perfect recovery up to 11 sparse entries, after which performance drops sharply.

To evaluate the practical impact of these two assumptions, we simulated data from the model in Eq. 3, with shared disturbances as in Eq. 9. We used 8 views, 6 common disturbances (two sub-Gaussians sampled from a generalized normal distribution with shape parameter $\beta = 2.5$; two super-Gaussians for $\beta = 1.5$; two Gaussians corresponding to $\beta = 2$), and 1000 samples. The scale matrices \mathbf{D}^i were drawn uniformly from the interval $[0.5, 2]$, the strictly lower triangular matrices \mathbf{T}^i were sampled from a Gaussian distribution and then permuted to form the causal matrices \mathbf{B}^i , and the variances Σ^i were sampled uniformly in $[0, 1]$. The experiment was repeated 30 times.

Fig. 9 reports the Spearman’s rank correlation between the true and estimated causal orderings as a function of the number of sparse entries in each \mathbf{T}^i , randomly chosen among the $\frac{6 \times 5}{2} = 15$ non-zero entries.

In the multiple (or view-specific) causal orderings setting (left panel), both ICA-LiMVAM variants exhibit a roughly linear drop in performance as sparsity increases. This behavior is expected, as Assumption 7 no longer holds once any entry in the strictly lower triangular part of \mathbf{T}^i is set to zero. PairwiseLiMVAM is not included in this setting, as it is designed under the assumption of a shared causal ordering.

In the shared causal ordering setting (right panel), all methods benefit from the shared structure and recover the true ordering \mathbf{P} reliably up to 11 sparse entries. Beyond this point, performance degrades sharply for all methods, reflecting the increasing violation of Assumption 3.

These results empirically support the practical relevance of Assumptions 3 and 7.

H.2. Real data experiments

H.2.1. DETAILS ON THE PREPROCESSING OF MEG DATA

The MEG data measures participants’ responses to auditory and visual stimuli. The auditory stimuli were binaural pure tones. The visual stimuli were checkerboards presented both to the left and right of a central fixation for 34-ms duration. This task leads to strong signal power modulations during the motor preparation and motor execution.

The original MEG data were acquired with 306 sensors, recorded at 1000 Hz, and band-pass filtered between 0.03 and 330 Hz. All MEG processing was done using the MNE-Python library (Gramfort et al., 2013; 2014) and we largely followed the pre-processing steps used in Power et al. (2023). We applied a Maxwell filter (Taulu and Simola, 2006) to improve data quality and a band-pass filter between 8 and 27 Hz to focus on power effects spread over the alpha and beta bands of the brain. This range of waves is supposed to be particularly active in sensorimotor tasks, especially during movement preparation and execution, and it typically shows a characteristic suppression (event-related desynchronization) during movement, followed by a rebound (event-related synchronization) after movement cessation, which is thought to reflect sensorimotor processing and inhibitory mechanisms. In particular, the suppression and rebound are reflected in the energies of the signals, not raw signals.

The data were then parsed into trials synchronized to each button press, with a duration of 4.5 s, including a 1.5 s pre-movement interval. The 4.5 s window length was selected to ensure a sufficient post-movement interval to capture the entire beta rebound response. Trials were excluded if the button press occurred more than 1 s after the audiovisual cue (indicating poor task performance) or if another button press occurred within the time window. Then, a baseline correction was applied using the pre-movement interval (-1.5 s, -1 s). The procedure led to about 60 trials per participant on average.

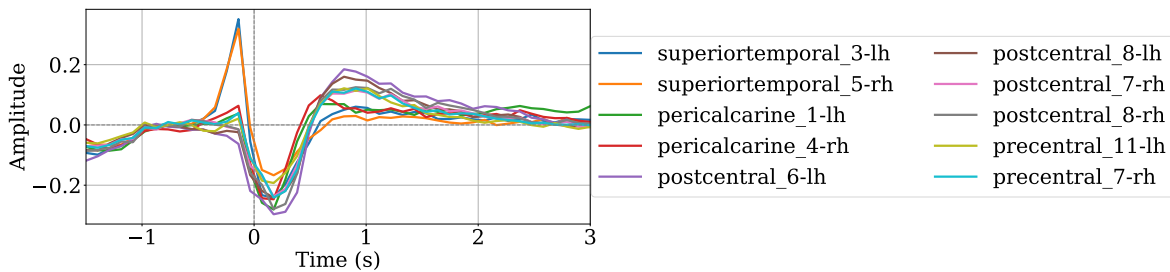


Figure 10. Example of time series obtained after the preprocessing. We averaged the data across subjects and trials, resulting in 10 time series, one for each brain region.

Next, we performed a cortical projection. First, each participant’s MRI (additionally given in the Cam-CAN database) was segmented using FreeSurfer (Dale et al., 1999). The segmentation provided a digitization of the cortical surface for source estimation, a transformation to the average brain (*i.e.* fsaverage) for spatial normalization and group statistics, and a boundary element model of the head to provide more accurate calculation of the forward solution. The inverse solution, based on the MNE method (Hämäläinen and Ilmoniemi, 1994), allowed to consider cortical region activations for further analysis.

We used the cortical parcellation from (Khan et al., 2018) to divide the cortical mantle (both hemispheres) into 448 distinct regions and summarized each region by an averaged time course. Then, we selected 10 of the 448 regions based on their known importance in sensorimotor tasks. Specifically, we picked for each hemisphere three regions in the motor cortex (two parcels in the “postcentral” and one in the “precentral”; visible in blue in Fig. 4a), one region in the auditory cortex (“superiortemporal” parcel; highlighted in pink), and one region in the visual cortex (“pericalcarine” parcel; not visible in the figure). Note again that the task for the participant is active as right index button presses are triggered by audiovisual stimulations.

For each participant, we thus extracted one time series for each of the 10 regions and fixed the number of trials to 40. Participants with less than 40 trials were discarded and extra trials were averaged. Participants for whom some of the parcels did not have any vertex in the source space were also discarded, resulting in a total of 98 available participants. We performed a Z-score normalization of the time series to correct the

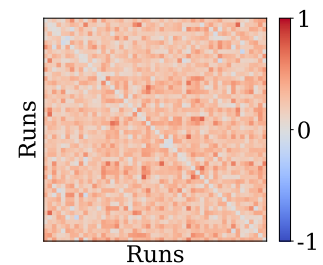


Figure 11. Pearson correlations between median causal effect matrices across 50 runs of ICA-LiMVAM-ML, each using a different random subset of 30 participants from the Cam-CAN cohort.

depth bias and, importantly, computed the Hilbert *envelope* of the signals to study modulations in cortical source power. Finally, the signals were centered and downsampled from 1000 Hz to 10 Hz due to the slowly changing nature of the envelopes (energies). The final dataset consisted of 98 subjects, 10 brain regions, and 1760 time points.

Figure 10 shows the typical time series obtained after preprocessing. For visualization, we averaged the time courses across participants and trials. As expected, we observe a prominent peak in the “superiortemporal” auditory regions shortly before the button press ($t = 0$ s), reflecting the stimulus onset. In the motor regions, a characteristic rebound pattern emerges: the signals decrease around the time of the button press, consistent with event-related desynchronization of beta rhythms, and subsequently increase after approximately 0.3 seconds, indicating beta resynchronization. This figure also labels the ten cortical parcels selected for analysis.

H.2.2. MEG EXPERIMENT USING ICA-LiMVAM-ML

In the following, we present additional analyses that extend the experiments described in Section 6.2.

Figure 12 displays median causal effect maps estimated using ICA-LiMVAM-ML on the Cam-CAN dataset. For each of the six panels, ICA-LiMVAM-ML was applied to a randomly chosen subset of 50% of the participants. These results are consistent with the patterns observed in Figure 4a, notably the frequent presence of directed connections between motor areas in opposite hemispheres. This further supports the hypothesis of consistent inter-hemispheric causal influences in sensorimotor processing.

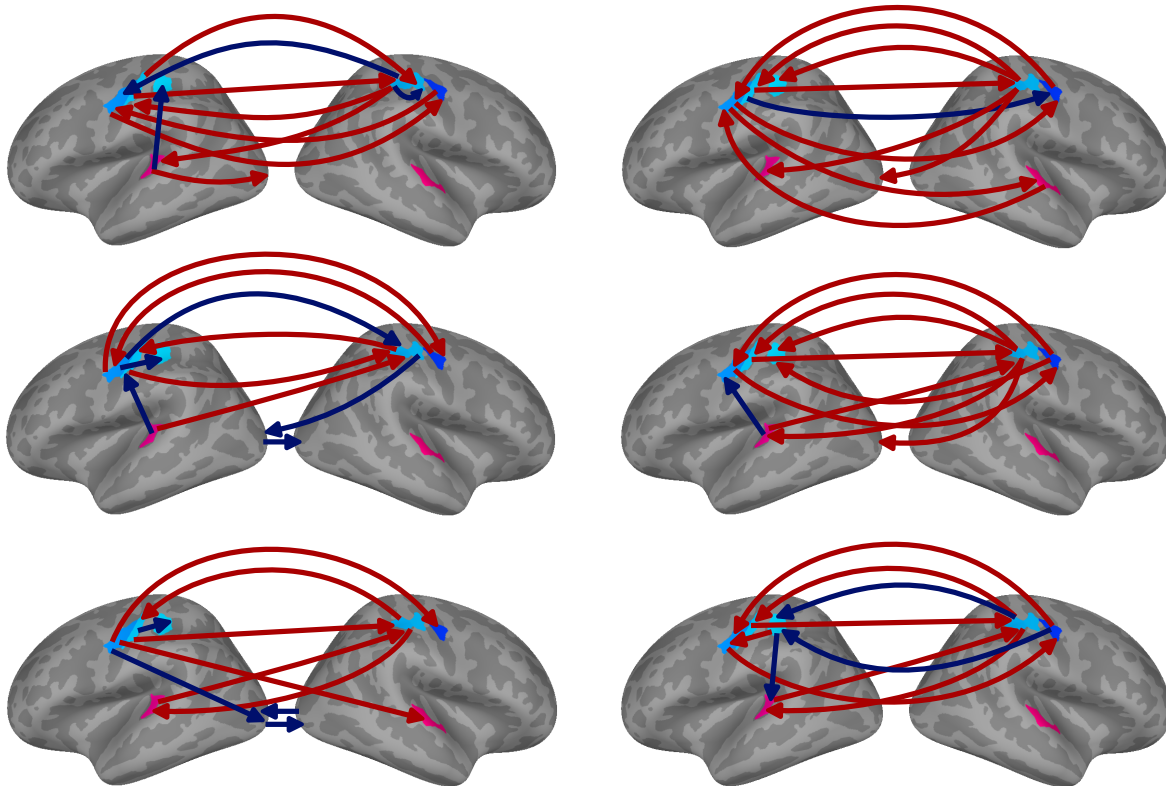


Figure 12. Top ten strongest median causal effects (estimated by ICA-LiMVAM-ML) of six different runs. Each run was performed on the data of 49 randomly chosen subjects. Red arrows represent positive effects and blue arrows negative effects.

Finally, we conducted the same robustness analysis for ICA-LiMVAM-ML as previously done for PairwiseLiMVAM (Figure 4b). Specifically, we ran ICA-LiMVAM-ML 50 times, each time using a randomly selected subset of 30 participants from the full cohort, and computed the Pearson correlations between the element-wise median of the estimated individual matrices B^i across runs. The resulting correlations are shown in Figure 11. While the correlations remain predominantly positive (with an average of 0.27), they are noticeably lower than those obtained with PairwiseLiMVAM. This suggests that, in terms of stability across subsets, PairwiseLiMVAM exhibits greater consistency than ICA-LiMVAM-ML.

I. Distinction between multi-view and multi-domain frameworks

Several frameworks study causal discovery from multiple related datasets that share the same causal structure. While *multi-view* methods (like ours) model *correlated* views arising from a joint (non-factorial) distribution, *multi-domain* methods treat datasets as *independent* domains drawn from related but distinct distributions. Consequently, multi-domain methods cannot leverage cross-view correlations and instead rely exclusively on distributional shifts across domains. *Multi-environment* methods are a special case of multi-domain methods, where distributional shifts arise from *interventions* (whether designed by the practitioner or arising from uncontrolled environment differences) on the causal mechanisms. In the following, we review some multi-domain/multi-environment methods that are related to our work.

Ghassami et al. (2018) propose a multi-domain method that, like ours, studies linear causal discovery from multiple datasets that share the same underlying DAG in the causally sufficient setting and without a non-Gaussianity assumption. However, they assume that causal mechanisms $\mathbb{P}(x_j^i | \mathbf{PA}_j^i)$ are *independent* (both within and across domains), where \mathbf{PA}_j^i denotes the parents of x_j^i . Our approach does not require such an independence assumption. For identifiability of the true DAG, their method further assumes that noise variances or causal weights vary across domains, and one of their two criteria additionally requires these changes to be sparse. In our paper, we show that such changes in noise variances are sufficient for identifiability (see Assumption 6), but not necessary when the views exhibit diverse correlations.

Adams et al. (2021) study multi-domain identifiability in the more complex setting of *latent confounders*. Their main contributions are necessary and sufficient conditions (“bottleneck” and “strong non-redundancies”) that are specific to the confounded setting and vanish in our causally sufficient case. In particular, in this causally sufficient case, their conditions reduce to assuming heterogeneous variances (see their Theorem 1), which is precisely one of the sufficient conditions in our Assumption 6. Moreover, they require the causal weights \mathbf{B}^i to *remain constant* across domains, which is a strong assumption we do not make. Thus, in the causally sufficient case, our method encompasses theirs and goes much further by leveraging correlations (across views and within a view), while allowing causal weights to differ. They suggested, in fact, that their assumptions were too strong in Section 3: “Note that this theorem gives sufficient conditions; our empirical results suggest that they are not necessary.”

Peters et al. (2016) propose a multi-environment method that aims to identify the parents of a *single target variable* rather than recovering the full graph, while assuming purely *Gaussian* distributions. Their approach relies on *invariance* of the target’s causal mechanism across environments, which implicitly requires knowing where interventions occurred (and that the target was not intervened upon), an assumption that is often unrealistic in practice.

Similarly, Perry et al. (2022) propose another multi-environment method that leverages invariance and can be more easily applied to causal discovery of the whole structure. However, their method considers the more complex case of *nonlinear* relationships. Like Ghassami et al. (2018), they rely on changes in $\mathbb{P}(x_j^i | \mathbf{PA}_j^i)$ across environments, but instead of measuring independence between mechanisms, they count how often these mechanisms change and select the DAG that minimizes this number. Their identifiability results require sparsity in the number of changes; in particular, not all causal mechanisms are allowed to differ across environments. In contrast, our theory accommodates arbitrarily many changes in \mathbf{B}^i and in the distribution of \mathbf{e}^i ; moreover, these changes are not required for identifiability if cross-view correlations are diverse.

We finally briefly mention two further works that are related to the multi-domain setting. A general framework, allowing for various kinds of interventions, was proposed by Mooij et al. (2020), but this was not claimed to improve identifiability. Sturma et al. (2023) proposed a related “multi-context” approach with the goal of estimating a representation as well as in causal representation learning (CRL).

In neuroimaging, this distinction between the two frameworks makes multi-view methods particularly well-suited to experiments centered around stimulus onsets (“phase-locked”, as in our MEG and fMRI experiments), while multi-domain methods are more naturally aligned with resting-state experiments (as in the fMRI experiment of Ghassami et al. (2018)).