

---

# An efficient search-and-score algorithm for ancestral graphs using multivariate information scores

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose a greedy search-and-score algorithm for ancestral graphs, which in-  
2 clude directed as well as bidirected edges, originating from unobserved latent  
3 variables. The normalized likelihood score of ancestral graphs is estimated in terms  
4 of multivariate information over relevant subsets of vertices,  $C$ , that are connected  
5 through collider paths confined to the ancestor set of  $C$ . For computational effi-  
6 ciency, the proposed two-step algorithm relies on local information scores limited  
7 to the close surrounding vertices of each node (step 1) and edge (step 2). This  
8 computational strategy is shown to outperform state-of-the-art causal discovery  
9 methods on challenging benchmark datasets.

## 10 1 Introduction

11 The likelihood function plays a central role in the selection of a graphical model  $\mathcal{G}$  based on  
12 observational data  $\mathcal{D}$ . Given  $N$  independent samples from  $\mathcal{D}$ , the likelihood  $\mathcal{L}_{\mathcal{D}|\mathcal{G}}$  that they might  
13 have been generated by the graphical model  $\mathcal{G}$  is given by [1],

$$\mathcal{L}_{\mathcal{D}|\mathcal{G}} = \frac{1}{Z_{\mathcal{D},\mathcal{G}}} \exp(-NH(p, q)) \quad (1)$$

14 where  $H(p, q) = -\sum_{\mathbf{x}} p(\mathbf{x}) \log q(\mathbf{x})$  is the cross-entropy between the empirical probability distribu-  
15 tion  $p(\mathbf{x})$  of the observed data  $\mathcal{D}$  and the theoretical probability distribution  $q(\mathbf{x})$  of the model  $\mathcal{G}$  and  
16  $Z_{\mathcal{D},\mathcal{G}}$  a data- and model-dependent factor ensuring proper normalization condition for finite dataset. In  
17 short, Eq.1 results from the asymptotic probability that the  $N$  independent samples,  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$ ,  
18 are drawn from the model distribution,  $q(\mathbf{x})$ , *i.e.*  $\mathcal{L}_{\mathcal{D}|\mathcal{G}} \equiv q(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) = \prod_i q(\mathbf{x}^{(i)})$ , rather  
19 than the empirical distribution,  $p(\mathbf{x})$ . This leads to,  $\log \mathcal{L}_{\mathcal{D}|\mathcal{G}} = \sum_i \log q(\mathbf{x}^{(i)})$ , which converges  
20 towards  $N \sum_{\mathbf{x}} p(\mathbf{x}) \log q(\mathbf{x}) = -NH(p, q)$  in the large sample size limit,  $N \rightarrow \infty$ , with  
21  $\log Z_{\mathcal{D},\mathcal{G}} = \mathcal{O}(\log N)$ .

22 The structural constraints of the model  $\mathcal{G}$  translate into the factorization form of the theoretical  
23 probability distribution,  $q(\mathbf{x})$  [2–6]. In particular, the probability distribution of Bayesian networks  
24 (BN) factorizes in terms of conditional probabilities of each variable given its parents, as  $q_{\text{BN}}(\mathbf{x}) =$   
25  $\prod_i q(x_i|\mathbf{pa}_{X_i})$ , where  $\mathbf{pa}_{X_i}$  denote the values of the parents of node  $X_i$  in  $\mathcal{G}$ ,  $\mathbf{Pa}_{X_i}$ . For Bayesian  
26 networks, the factors of the model distribution,  $q(x_i|\mathbf{pa}_{X_i})$ , can be directly estimated with the  
27 empirical conditional probabilities of each node given its parents as,  $q(x_i|\mathbf{pa}_{X_i}) \equiv p(x_i|\mathbf{pa}_{X_i})$ ,  
28 leading to the well known estimation of the likelihood function in terms of conditional entropies  
29  $H(X_i|\mathbf{Pa}_{X_i}) = -\sum_{\mathbf{x}} p(x_i, \mathbf{pa}_{X_i}) \log p(x_i|\mathbf{pa}_{X_i})$ ,

$$\mathcal{L}_{\mathcal{D}|\mathcal{G}_{\text{BN}}} = \frac{1}{Z_{\mathcal{D},\mathcal{G}_{\text{BN}}}} \exp\left(-N \sum_{X_i \in \mathcal{V}}^{\text{vertices}} H(X_i|\mathbf{Pa}_{X_i})\right) \quad (2)$$

30 This paper concerns the experimental setting for which some variables of the underlying Bayesian  
 31 model are not observed. This frequently occurs in practice for many applications. We derive an  
 32 explicit likelihood function for the class of ancestral graphs, which include directed as well as  
 33 bidirected edges, arising from the presence of unobserved latent variables. Tian and Pearl 2002 [7]  
 34 showed that the probability distribution of such graphs factorizes into c-components including subsets  
 35 of variables connected through bidirected paths (*i.e.* containing only bidirected edges). Richardson  
 36 2009 [6] later proposed a refined factorization of the model distribution of the broader class of acyclic  
 37 directed mixed graphs in terms of conditional probabilities over “head” and “tail” subsets of variables  
 38 within each ancestrally closed subsets of vertices. However, unlike with Bayesian networks, the  
 39 contributions of c-components or head-and-tail factors to the likelihood function cannot simply be  
 40 estimated in terms of empirical distribution  $p(\mathbf{x})$ , as shown below. This leaves the likelihood function  
 41 of ancestral graphs difficult to estimate from empirical data, in general, although iterative methods  
 42 have been developed when the data is normally distributed [8–13].

43 The present paper provides an explicit decomposition of the likelihood function of ancestral graphs  
 44 in terms of multivariate cross-information over relevant ‘*ac*-connected’ subsets of variables, Figs. 1.,  
 45 which do not rely on the head-and-tail factorization but coincide with the parametrizing sets [14]  
 46 derived from the head-and-tail factorization. It suggests a natural estimation of these relevant  
 47 contributions to the likelihood function in terms of empirical distribution  $p(\mathbf{x})$ . This result extends  
 48 the likelihood expression of Bayesian Networks (Eq. 2) to include the effect of unobserved latent  
 49 variables and enables the implementation of a greedy search-and-score algorithm for ancestral graphs.  
 50 For computational efficiency, the proposed two-step algorithm relies on local information scores  
 51 limited to the close surrounding vertices of each node (step 1) and edge (step 2). This computational  
 52 strategy is shown to outperform state-of-the-art causal discovery methods on challenging benchmark  
 53 datasets.

## 54 2 Theoretical results

### 55 2.1 Multivariate cross-entropy and cross-information

56 The theoretical result of the paper (Theorem 1) is expressed in terms of multivariate cross-information  
 57 derived from multivariate cross-entropies through the Inclusion-Exclusion Principle. The same  
 58 expressions can be written between multivariate information and multivariate entropies by simply  
 59 substituting  $q(\{x_i\})$  with  $p(\{x_i\})$  in the equations below and will be used to estimate the likelihood  
 60 function of ancestral graphs (Proposition 3).

61 As recalled above, the cross-entropy between  $m$  variables,  $\mathbf{V} = \{X_1, \dots, X_m\}$ , is defined as,

$$H(\mathbf{V}) = - \sum_{\{x_i\}} p(x_1, \dots, x_m) \log q(x_1, \dots, x_m) \quad (3)$$

62 where  $p(\{x_i\})$  is the empirical joint probability distribution of the variables  $\{X_i\}$  and  $q(\{x_i\})$  the  
 63 joint probability distribution of the model. Bayes formula,  $q(\{x_i\}, \{y_j\}) = q(\{x_i\}|\{y_j\}) q(\{y_j\})$ ,  
 64 directly translates into the definition of conditional cross-entropy through the decomposition,  
 65

$$H(\{X_i\}, \{Y_j\}) = H(\{X_i\}|\{Y_j\}) + H(\{Y_j\}) \quad (4)$$

66 Multivariate (cross) information,  $I(\mathbf{V}) \equiv I(X_1; \dots; X_m)$ , are defined from multivariate (cross)  
 67 entropies through Inclusion-Exclusion formulas over all subsets of variables [15–18] as,

$$\begin{aligned} I(X) &= H(X) \\ I(X; Y) &= H(X) + H(Y) - H(X, Y) \\ I(X; Y; Z) &= H(X) + H(Y) + H(Z) - H(X, Y) - H(X, Z) - H(Y, Z) + H(X, Y, Z) \\ I(\mathbf{V}) &= - \sum_{S \subseteq \mathbf{V}} (-1)^{|S|} H(S) \end{aligned} \quad (5)$$

68 where the semicolon separators are needed to distinguish multipoint (cross) information from joint  
 69 variables as  $\{X, Z\}$  in  $I(\{X, Z\}; Y) = I(X; Y) + I(Z; Y) - I(X; Y; Z)$ . Below, implicit separators  
 70 between non-conditioning variables in multivariate (cross) information will always correspond to  
 71 semicolons, *e.g.* as in  $I(\mathbf{V})$  in Eq. 5. Unlike multivariate (cross) entropies, which are always positive,

72  $H(X_1, \dots, X_k) \geq 0$ , multivariate (cross) information,  $I(X_1; \dots; X_k)$ , can be positive or negative  
 73 for  $k \geq 3$ , while they remain always positive for  $k < 3$ , i.e.  $I(X; Y) \geq 0$  and  $I(X) \geq 0$ .

74 In turn, multivariate (cross) entropies can be expressed through the Principle of Inclusion-Exclusion  
 75 into the same expression form but in terms of multivariate (cross) information,

$$H(\mathbf{V}) = - \sum_{\mathcal{S} \subseteq \mathbf{V}} (-1)^{|\mathcal{S}|} I(\mathcal{S}), \quad (6)$$

76 Conditional multivariate (cross) information  $I(\mathbf{V}|Z)$  are defined similarly as multivariate (cross)  
 77 information  $I(\mathbf{V})$  but in terms of conditional (cross) entropies as,

$$I(\mathbf{V}|Z) = - \sum_{\mathcal{S} \subseteq \mathbf{V}} (-1)^{|\mathcal{S}|} H(\mathcal{S}|Z) \quad (7)$$

78 Eqs. 5 & 7 lead to a decomposition rule relative to a variable  $Z$ , Eq. 8, which can be conditioned  
 79 on a set of joint variables,  $\mathbf{A} = \{A_1, \dots, A_m\}$ , with implicit comma separators for conditioning  
 80 variables in Eq. 9,

$$I(\mathbf{V}) = I(\mathbf{V}|Z) + I(\mathbf{V}; Z) \quad (8)$$

$$I(\mathbf{V}|\mathbf{A}) = I(\mathbf{V}|Z, \mathbf{A}) + I(\mathbf{V}; Z|\mathbf{A}) \quad (9)$$

81 Alternatively, conditional (cross) information, such as  $I(X; Y|\mathbf{A})$ , can be expressed in terms of  
 82 non-conditional (cross) entropies using Eq. 4,

$$\begin{aligned} I(X; Y|\mathbf{A}) &= H(X|\mathbf{A}) + H(Y|\mathbf{A}) - H(X, Y|\mathbf{A}) \\ &= H(X, \mathbf{A}) + H(Y, \mathbf{A}) - H(X, Y, \mathbf{A}) - H(\mathbf{A}) \end{aligned} \quad (10)$$

83 which can in turn be expressed in terms of non-conditional (cross) information as,

$$\begin{aligned} I(X; Y|\mathbf{A}) &= I(X; Y) - \dots (-1)^k \sum_{i_1 < \dots < i_k} I(X; Y; A_{i_1}; \dots; A_{i_k}) + \dots (-1)^m I(X; Y; A_1; \dots; A_m) \\ &= \sum_{\mathcal{S}' \subseteq \mathcal{S}}^{X, Y \in \mathcal{S}'} (-1)^{|\mathcal{S}'|} I(\mathcal{S}'), \end{aligned} \quad (11)$$

84 where  $\mathcal{S} = \{X, Y\} \cup \mathbf{A}$ . This corresponds, up to an opposite sign, to *all (cross) information terms*  
 85 *including both  $X$  and  $Y$*  in the expression of the multivariate (cross) entropy,  $H(X, Y, \mathbf{A})$ , Eq. 6.

## 86 2.2 Graphs and connection criteria

### 87 2.2.1 Directed mixed graphs and ancestral graphs

88 Two vertices are said to be **adjacent** if there is an edge (of any type) between them,  $X * \rightarrow Y$ , where  
 89  $*$  stands for any (head or tail) end mark.  $X$  and  $Y$  are said to be **neighbors** if  $X - Y$ , **parent** and  
 90 **child** if  $X \rightarrow Y$  and **spouses** if  $X \longleftrightarrow Y$  in  $\mathcal{G}$ .

91 A **path** in  $\mathcal{G}$  is a sequence of distinct vertices  $V_1, \dots, V_n$  consecutively adjacent in  $\mathcal{G}$ , as,  
 92  $V_1 * \rightarrow V_2 * \rightarrow \dots * \rightarrow V_{n-1} * \rightarrow V_n$ . In particular, a **collider path** between  $V_1$  and  $V_n$  has the form  
 93  $V_1 * \rightarrow V_2 \leftarrow \dots \leftarrow V_{n-1} \leftarrow * V_n$  and a **directed path** corresponds to  $V_1 \rightarrow V_2 \rightarrow \dots \rightarrow V_n$ .

94  $X$  is called an **ancestor** of  $Y$  and  $Y$  a **descendant** of  $X$  if  $X = Y$  or there is a **directed path** from  
 95  $X$  to  $Y$ ,  $X \rightarrow \dots \rightarrow Y$ .  $\text{An}_{\mathcal{G}}(Y)$  denotes the **set of ancestors** of  $Y$  in  $\mathcal{G}$ . By extension, for any  
 96 subset of vertices,  $\mathcal{C} \subseteq \mathbf{V}$ ,  $\text{An}_{\mathcal{G}}(\mathcal{C})$  denotes the set of ancestors for all  $Y \in \mathcal{C}$  in  $\mathcal{G}$ .

97 A **directed mixed graph** is a vertex-edge graph  $\mathcal{G} = (\mathbf{V}, \mathbf{E})$  that can contain two types of edges:  
 98 directed ( $\rightarrow$ ) and bidirected ( $\longleftrightarrow$ ) edges.

99 A **directed cycle** occurs in  $\mathcal{G}$  when  $X \in \text{An}_{\mathcal{G}}(Y)$  and  $X \leftarrow Y$ . An **almost directed cycle** occurs  
 100 when  $X \in \text{An}_{\mathcal{G}}(Y)$  and  $X \longleftrightarrow Y$ .

101 **Definition 1.** An **ancestral graph** is a directed mixed graph:

- 102 i) without directed cycles;
- 103 ii) without almost directed cycles.

104 An **ancestral graph** is said to be **maximal** if every missing edge corresponds to a structural indepen-  
 105 dence. If an ancestral graph  $\mathcal{G}$  is not maximal, there exists a unique maximal ancestral graph  $\tilde{\mathcal{G}}$  by  
 106 adding bidirected edges to  $\mathcal{G}$  [8].

107 **2.2.2 *ac*-connecting paths and *ac*-connected subsets**

108 Let us now define **ancestor collider connecting paths** or ***ac*-connecting paths**, which entail simpler  
 109 path connecting criterion than the traditional **m-connecting criterion**, discussed in the Appendix A.  
 110 Yet, ***ac*-connecting paths** and ***ac*-connected subsets** will turn out to be directly relevant to character-  
 111 ize the likelihood decomposition and Markov equivalent classes of ancestral graphs.

112 **Definition 2.** [*ac*-connecting path] An *ac*-connecting path between  $X$  and  $Y$  given a subset of  
 113 variables  $C$  (possibly including  $X$  and  $Y$ ) is a collider path,  $X * \rightarrow Z_1 \leftarrow \dots \leftarrow Z_K \leftarrow * Y$ ,  
 114 with all  $Z_i \in \text{An}_{\mathcal{G}}(\{X, Y\} \cup C)$ , that is, with  $Z_i$  in  $C$  or connected to  $\{X, Y\} \cup C$  by an ancestor  
 115 path, *i.e.*  $Z_i \rightarrow \dots \rightarrow T$  with  $T \in \{X, Y\} \cup C$ .

116 **Definition 3.** [*ac*-connected subset] A subset  $C$  is said to be *ac*-connected if  $\forall X, Y \in C$ ,  $X$  and  
 117  $Y$  are connected (through any type of edge) or there is an *ac*-connecting path between  $X$  and  $Y$   
 118 given  $C$ .

119 **2.3 Likelihood decomposition of ancestral graphs**

120 **Theorem 1.** [**likelihood of ancestral graphs**] *The cross-entropy  $H(p, q)$  and likelihood  $\mathcal{L}_{\mathcal{D}|\mathcal{G}}$  of an*  
 121 *ancestral graph  $\mathcal{G}$  is decomposable in terms of multivariate cross-information,  $I(C)$ , summed over*  
 122 *all *ac*-connected subsets of variables,  $C$  (Definition 3),*

$$\begin{aligned}
 H(p, q) &= - \sum_{C \subseteq V}^{\text{ac-connected}} (-1)^{|C|} I(C) \\
 \mathcal{L}_{\mathcal{D}|\mathcal{G}} &= \frac{1}{Z_{\mathcal{D}, \mathcal{G}}} \exp \left( N \sum_{C \subseteq V}^{\text{ac-connected}} (-1)^{|C|} I(C) \right) \quad (12)
 \end{aligned}$$

123 *where  $N$  is the number of iid samples in the dataset  $\mathcal{D}$  and  $Z_{\mathcal{D}, \mathcal{G}}$  a data- and model-dependent*  
 124 *normalization constant.*

125 The proof of Theorem 1 is left to the Appendix B. It is based on a partition of the cross-entropy (Eq. 6)  
 126 into cross-information contributions from *ac*-connected and non-*ac*-connected subsets of variables,  
 127 which do not rely on head-and-tail factorizations. Hu and Evans [14] proposed an equivalent result  
 128 (Proposition 3.3 in [14]) with a proof using head-and-tail decomposition to define parametrizing  
 129 sets, which happen to coincide with the *ac*-connected sets defined here (Definition 3). Theorem 1  
 130 characterizes in particular the Markov equivalence class of ancestral graphs [8, 19–24] as,

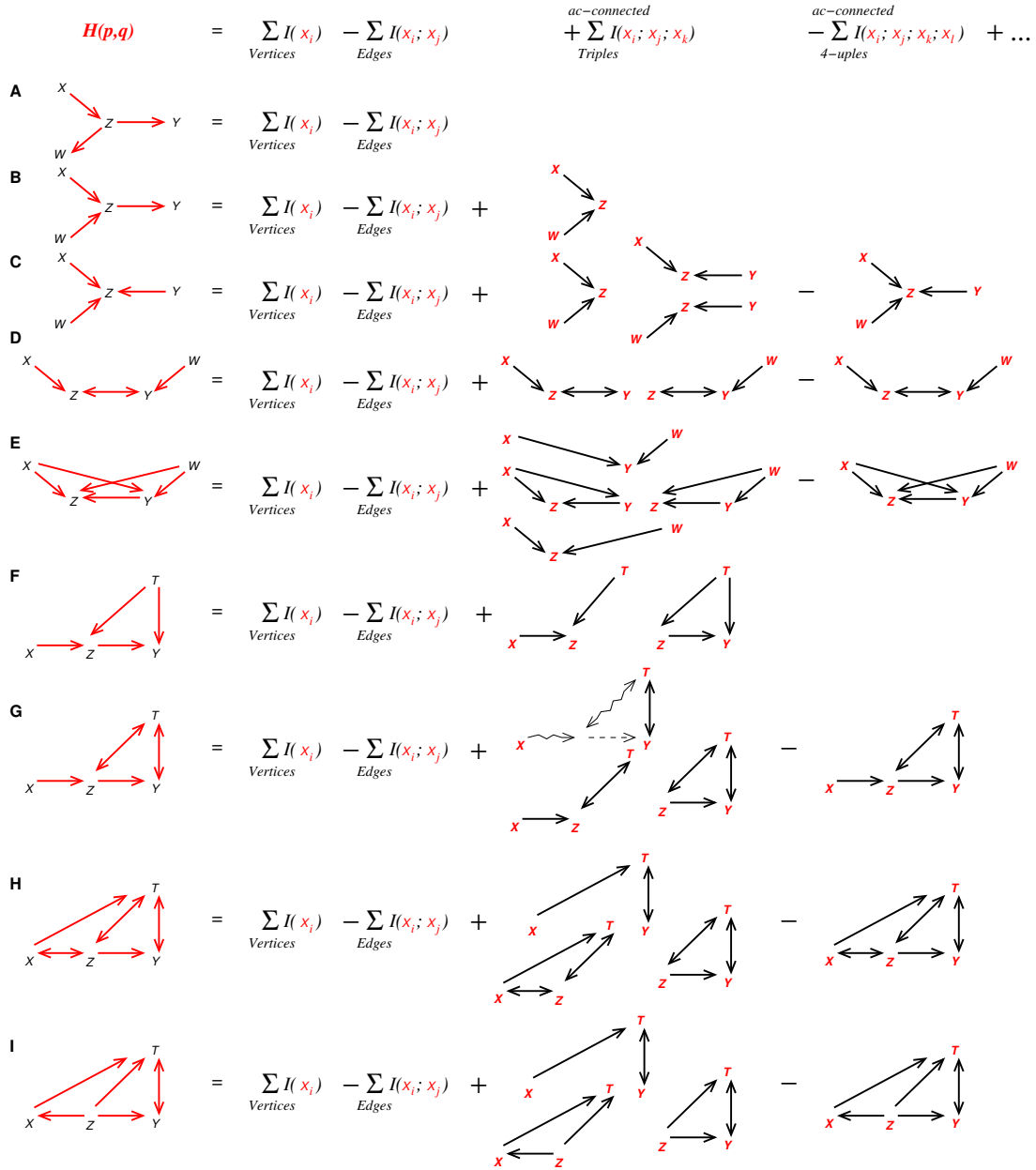
131 **Corollary 2.** *Two ancestral graphs are Markov equivalent if and only if they have the same *ac*-*  
 132 *connected subsets of vertices.*

133 Note, in particular, that Eq. 12 holds for *maximal ancestral graphs* (MAG), for which all pairs of  
 134 *ac*-connected variables are connected by an edge, and their Markov equivalent representatives, the  
 135 *partial ancestral graphs* (PAG) [8, 25–27].

136 **Proposition 3.** The likelihood decomposition of ancestral graphs (Eq. 12, Theorem 1) can be  
 137 estimated by replacing the model distribution  $q$  by the empirical distribution  $p$  in the retained  
 138 multivariate cross-information terms  $I(C)$  corresponding to all *ac*-connected subsets of variables,  $C$ .

139 Hence, Proposition 3 amounts to estimating all relevant cross-information terms in the likelihood  
 140 function with the corresponding multivariate information terms computed from the available data,  
 141 while assuming by construction that the model distribution obeys all local and global conditional inde-  
 142 pendences entailed by the ancestral graph. The corresponding factorization of the model distribution  
 143 can be expressed in terms of empirical distribution, assuming positive distributions, see Appendix C.

144 Fig. 1 illustrates the cross-entropy decomposition for a few graphical models in terms of cross-  
 145 information contributions from their *ac*-connected subsets of vertices. In particular, an unshielded  
 146 non-collider (*e.g.*  $X \rightarrow Z \rightarrow W$ , Fig. 1A), is less likely (*i.e.* higher cross-entropy) than an unshielded  
 147 collider or ‘v-structure’ (*e.g.*  $X \rightarrow Z \leftarrow W$ , Fig. 1B), if the corresponding three-point information  
 148 term is negative,  $I(X; Z; W) < 0$ , in agreement with earlier results [28, 29]. However, this early  
 149 approach, exploiting the sign and magnitude of three-point information to orient v-structures, does  
 150 not include higher order terms involving multiple v-structures, which can lead to orientation conflicts  
 151 between unshielded triples, in practice. Resolving such orientation conflicts requires to include



**Figure 1: Cross-entropy decomposition of ancestral graphs.** Examples of cross-entropy decomposition of ancestral graphs (red edges, lhs) in terms of relevant multivariate cross-information contributions  $I(C)$  with  $C \subseteq V$  (red nodes, rhs). Simple graphs: **(A)** without unshielded colliders, **(B)** with a single or non-overlapping unshielded colliders, **(C)** with overlapping unshielded colliders through three or more (conditionally) independent parents or **(D)** through a two-(or more)-collider path. **(E)** Bayesian graph corresponding to the head-and-tail factorization of the two-collider path in **(D)** estimated using the empirical distribution  $p(\cdot)$ , see Appendix C. **(F)** Simple Bayesian graph not Markov equivalent to an ancestral graph **(G)** sharing the same edges and unshielded collider [24]. Solid black edges correspond to direct connections or collider paths confined to the corresponding  $ac$ -connected subset  $C$ , while wiggly edges indicate collider paths extending beyond  $C$  yet indirectly connected to  $C$  by an ancestor path, marked with dashed edges, see Definition 2. By contrast, graphs **H** and **I** illustrate the fact that collider paths may not be unique nor conserved between two Markov equivalent graphs (*i.e.* sharing the same cross-information terms) [24].

152 information contributions from higher-order  $ac$ -connected subgraphs, such as star-like  $ac$ -connected  
 153 subsets including three or more parents, Fig. 1C. Similarly, the cross-entropies of collider paths  
 154 involving several colliders also include higher-order terms, as with the simple example of a two-

155 collider path, Fig. 1D. By contrast, the cross-entropy based on the head-and-tail factorization of the  
 156 same two-collider path, *i.e.*  $q(x, z, y, w) = q(z, y|x, w)q(x)q(w)$  [6], is found to be equivalent to  
 157 the cross-entropy of a Bayesian graph without bidirected edge, Fig. 1E, when estimated with the  
 158 empirical distribution  $p(\cdot)$ , see Appendix C. This observation illustrates the difficulty to estimate the  
 159 likelihood functions of ancestral graphs using head-and-tail factorization.

160 Further examples of graphical models, Figs. 1F-I, show the relative simplicity of the decomposition  
 161 with only few (non-trivial) *ac*-connected contributing subsets  $\mathcal{C}$  with  $|\mathcal{C}| \geq 3$ , as compared to  
 162 the much larger number of non-*ac*-connected non-contributing subsets, that cancel each other by  
 163 construction due to conditional independence constraints of the underlying model. Note, in particular,  
 164 that most contributing multivariate information  $I(\mathcal{C})$  only concern direct connections or collider  
 165 paths within a single component subgraph induced by  $\mathcal{C}$  (solid line edges in Fig. 1). However,  
 166 occasionally, collider paths extending beyond  $\mathcal{C}$  into  $\text{An}_{\mathcal{G}}(\mathcal{C}) \setminus \mathcal{C}$  (marked with wiggly edges) with  
 167 corresponding ancestor path(s) (marked with dashed edges) do occur, as shown in Fig. 1G.

168 In addition, the present information-theoretic decomposition of the likelihood of ancestral graphs  
 169 can readily distinguish their Markov equivalence classes according to Corollary 2. For instance, the  
 170 ancestral graphs of Fig. 1F and Fig. 1G, despite sharing the same edges and the same unshielded  
 171 collider ( $X \rightarrow Z \leftarrow T$ ), turn out not to be Markov equivalent, as discussed in [24]. Indeed, their  
 172 cross-entropy decompositions differ by two *ac*-connected contributing terms: a three-point cross  
 173 information  $I(X; Y; T)$  with a collider path not confined in  $\mathcal{C}$  (*i.e.*  $X \rightsquigarrow Z \leftrightarrow T \leftrightarrow Y$  and  
 174 corresponding ancestor path  $Z \dashrightarrow Y$ ) and a four-point information term  $I(X; Y; Z; T)$  due to  
 175 the two-collider path ( $X \rightarrow Z \leftrightarrow T \leftrightarrow Y$ ). More quantitatively, it shows that the graph of  
 176 Fig. 1G with a two-collider path is more likely than the graph of Fig. 1F whenever  $I(X; Y; T) -$   
 177  $I(X; Y; Z; T) = I(X; Y; T|Z) = I(X; Y|Z) - I(X; Y|Z, T) < 0$ . Finally, the Markov equivalent  
 178 graphs of Fig. 1H and Fig. 1I, also due to [24], illustrate the fact that the actual ancestor collider path  
 179 between unconnected pairs does not need to be unique nor conserved between Markov equivalent  
 180 graphs (as long as their cross-entropies share the same multivariate cross-information decomposition).

### 181 3 Efficient search-and-score causal discovery using local information scores

182 The likelihood estimation of ancestral graphs (Theorem 1 and Proposition 3) enables the implemen-  
 183 tation of a search-and-score algorithm for this broad class of graphs, which has attracted a number  
 184 of contributions recently [11–13, 30–32]. Our specific objective is not to develop an exact method  
 185 limited to simple graphical models with a few nodes and small datasets but to implement an efficient  
 186 and reliable heuristic method applicable to more challenging graphical models and large datasets.

187 Indeed, search-and-score structure learning methods need to rely on heuristic rather than exhaustive  
 188 search, in general, given that the number of ancestral graphs grows super-exponentially as the number  
 189 of vertices increases. This can be implemented for instance with a Monte Carlo algorithmic scheme  
 190 with random restarts, which efficiently probes relevant graphical models. Here, we opt, instead, to  
 191 use the prediction of an efficient hybrid causal discovery method, MIIC [29, 33, 34], as starting point  
 192 for a subsequent search-and-score approach based on the proposed likelihood estimation of ancestral  
 193 graphs (Eq. 12 and Proposition 3).

194 Moreover, while the likelihood decomposition of ancestral graphs may involve extended *ac*-connected  
 195 subsets of variables, as illustrated in Fig. 1, we aim to implement a computationally efficient search-  
 196 and-score causal discovery method based on approximate local scores limited to the close surrounding  
 197 vertices of each node and edge. Yet, while MIIC only relies on unshielded triple scores, the novel  
 198 search-and-score extension, MIIC\_search&score, uses also higher-order local information scores to  
 199 compare alternative subgraphs, as detailed below.

200 The proposed method is shown to outperform MIIC and other state-of-the-art causal discovery  
 201 methods on challenging datasets including latent variables.

#### 202 3.1 MIIC, an hybrid causal discovery method based on unshielded triple scores

203 MIIC is an hybrid causal discovery method combining constraint-based and information-theoretic  
 204 frameworks [29, 35]. Unlike traditional constraint-based methods [4, 5], MIIC does not directly  
 205 attempt to uncover conditional independences but, instead, iteratively substracts the most significant  
 206 three-point (conditional) information contributions of successive contributors,  $A_1, A_2, \dots, A_n$ , from

207 the mutual information between each pair of variables,  $I(X; Y)$ , as,

$$I(X; Y) - I(X; Y; A_1) - I(X; Y; A_2|A_1) - \dots - I(X; Y; A_n|\{A_i\}_{n-1}) = I(X; Y|\{A_i\}_n) \quad (13)$$

208 where  $I(X; Y; A_k|\{A_i\}_{k-1}) > 0$  is the *positive* information contribution from  $A_k$  to  $I(X; Y)$   
 209 [28, 36]. Conditional independence is eventually established when the residual conditional mutual  
 210 information on the right hand side of Eq. 13,  $I(X; Y|\{A_i\}_n)$ , becomes smaller than a complexity  
 211 term, *i.e.*  $k_{X;Y|\{A_i\}}(N) \geq I(X; Y|\{A_i\}_n) \geq 0$ , which depends on the considered variables and  
 212 sample size  $N$ .

213 This leads to an undirected skeleton, which MIIC then (partially) orients based on the sign and  
 214 amplitude of the regularized conditional 3-point information terms [28, 29]. In particular, negative  
 215 conditional 3-point information terms,  $I(X; Y; Z|\{A_i\}) < 0$ , correspond to the signature of causality  
 216 in observational data [28] and lead to the prediction of a v-structure,  $X \rightarrow Z \leftarrow Y$ , if  $X$  and  $Y$   
 217 are not connected in the skeleton. By contrast, a positive conditional 3-point information term,  
 218  $I(X; Y; Z|\{A_i\}) > 0$ , implies the absence of a v-structure and suggests to propagate the orientation  
 219 of a previously directed edge  $X \rightarrow Z - Y$  as  $X \rightarrow Z \rightarrow Y$ .

220 In practice, MIIC’s strategy to circumvent spurious conditional independences significantly improves  
 221 recall, that is, the fraction of correctly recovered edges, compared to traditional constraint-based  
 222 methods [28, 29]. Yet, MIIC only relies on unshielded triple scores to reliably uncover significant  
 223 contributors and orient v-structures, as outlined above. MIIC has been recently improved to ensure  
 224 the consistency of the separating set in terms of indirect paths in the final skeleton or (partially)  
 225 oriented graphs [37, 34] and to improve the reliability of predicted orientations [33, 34].

226 The predictions of this recent version of MIIC, which include three type of edges (directed, bidirected  
 227 and undirected), have been used as starting point for the subsequent local search-and-score method  
 228 implemented in the present paper.

## 229 3.2 New search-and-score method based on higher-order local information scores

230 Starting from the structure predicted by MIIC, as detailed above, MIIC\_search&score method  
 231 proceeds in two steps.

### 232 3.2.1 Step 1: Node scores for edge orientation priming and edge removal

233 The first step consists in minimizing a node score corresponding to the local normalized log likelihood  
 234 of each node w.r.t. its possible parents or spouses amongst the connected nodes predicted by MIIC.  
 235 To this end, the node score assesses the conditional entropy of each node w.r.t. a selection of  
 236 parents, spouses or neighbors,  $\mathbf{Pa}'_{X_i} \subseteq \mathbf{Pa}_{X_i} \cup \mathbf{Sp}_{X_i} \cup \mathbf{Ne}_{X_i}$ , and a factorized Normalized Maximum  
 237 Likelihood (fNML) regularization [28], see Appendix D for details,

$$\text{Score}_n(X_i) = H(X_i|\mathbf{Pa}'_{X_i}) + \frac{1}{N} \sum_j^{q_{x_i}} \log \mathcal{C}_{n_j}^{r_{x_i}} \quad (14)$$

238 where  $q_{x_i}$  corresponds to the combination of levels of  $\mathbf{Pa}'_{X_i}$ , while  $r_{x_i}$  is the number of levels of  $X_i$ ,  
 239 and  $n_j$  the number of samples corresponding to a particular combination of levels  $j$  in each summand,  
 240 with  $\sum_j n_j = N$ , the total number of samples.  $\log \mathcal{C}_{n_j}^{r_{x_i}}$  is the fNML regularization cost summed  
 241 over all combinations of levels,  $q_{x_i}$ , [38, 39], see Appendix D.

242 This first algorithm is looped over each node, priming the orientations of their surrounding edges (as  
 243 directed, bidirected or undirected), until convergence. Edges without orientation priming at either  
 244 extremity are removed at the end of Step 1.

### 245 3.2.2 Step 2: Edge orientation scores

246 The second step consists in minimizing an edge orientation score corresponding to the local normal-  
 247 ized log likelihood of each edge w.r.t. its nodes’ parents and spouses inferred in Step 1. To this end,  
 248 the edge score assesses the conditional information and a fNML complexity cost with respect to the  
 249 type of orientation, given three sets of parents and spouses of  $X$  and  $Y$ , *i.e.*  $\mathbf{Pa}'_{XY} = \mathbf{Pa}_X \cup \mathbf{Sp}_Y \setminus Y$ ,  
 250  $\mathbf{Pa}'_{YX} = \mathbf{Pa}_Y \cup \mathbf{Sp}_X \setminus X$  and  $\mathbf{Pa}'_{XY} = \mathbf{Pa}'_{XY} \cup \mathbf{Pa}'_{YX}$  with their corresponding combinations of

251 levels,  $q_{y \setminus x}$ ,  $q_{x \setminus y}$  and  $q_{x \setminus y}$ . These orientation scores, listed in Table 1, include symmetrized fNML com-  
 252 plexity terms to enforce Markov equivalence, if  $X$  and  $Y$  share the same parents or spouses (excluding  
 253  $X$  and  $Y$ ), see Appendix D. Indeed, all three scores become equals if  $\mathbf{Pa}'_{Y \setminus X} = \mathbf{Pa}'_{X \setminus Y} = \mathbf{Pa}'_{XY}$   
 254 implying also the same combinations of parent and spouse levels,  $q_{y \setminus x} = q_{x \setminus y} = q_{x \setminus y}$ .

255 This second algorithm is looped over each edge to compute an orientation score decrement, given the  
 256 orientations of its surrounding edges. The orientation change corresponding to the largest orientation  
 257 score decrement is then chosen at each iteration until convergence or until a limit cycle is reached  
 258 and stopped at the lowest sum of local orientation scores.

Table 1: Local scores for the orientation of a single directed or bidirected edge.

Edge	Information	Symmetrized fNML complexity (Markov equivalent)
$X \rightarrow Y$	$-I(X; Y   \mathbf{Pa}'_{Y \setminus X})$	$\frac{1}{2N} \left( \sum_j^{q_{x \setminus y} r_y} \log C_{n_j}^{r_x} - \sum_j^{q_{x \setminus y}} \log C_{n_j}^{r_x} + \sum_j^{q_{y \setminus x} r_x} \log C_{n_j}^{r_y} - \sum_j^{q_{y \setminus x}} \log C_{n_j}^{r_y} \right)$
$X \leftarrow Y$	$-I(X; Y   \mathbf{Pa}'_{X \setminus Y})$	$\frac{1}{2N} \left( \sum_j^{q_{x \setminus y} r_y} \log C_{n_j}^{r_x} - \sum_j^{q_{x \setminus y}} \log C_{n_j}^{r_x} + \sum_j^{q_{y \setminus x} r_x} \log C_{n_j}^{r_y} - \sum_j^{q_{y \setminus x}} \log C_{n_j}^{r_y} \right)$
$X \leftrightarrow Y$	$-I(X; Y   \mathbf{Pa}'_{XY})$	$\frac{1}{2N} \left( \sum_j^{q_{xy} r_y} \log C_{n_j}^{r_x} - \sum_j^{q_{xy}} \log C_{n_j}^{r_x} + \sum_j^{q_{yx} r_x} \log C_{n_j}^{r_y} - \sum_j^{q_{yx}} \log C_{n_j}^{r_y} \right)$

## 259 4 Experimental results

260 We first tested whether MIIC\_search&score orientation scores (Table 1) effectively predicts bidirected  
 261 orientations on three simple ancestral models, Fig. 3, when the end nodes do not share the same  
 262 parents (Fig. 3, Model 1), share some parents (Fig. 3, Model 2) or when the bidirected edge is part of  
 263 a longer than two-collider paths (Fig. 3, Model 3). The prediction of the edge orientation scores are  
 264 summarized in Table 3, Appendix E, and show good predictions for large enough datasets.

265 Beyond these simple examples, focussing on the discovery of bidirected edges in small toy models  
 266 of ancestral graphs, we also analyzed more challenging benchmarks from the bnlearn repository  
 267 [40], Fig. 2. They concern ancestral graphs obtained by hiding up to 20% of variables in Bayesian  
 268 Networks of increasing complexity (number of nodes and parameters), such as Alarm (37 nodes, 46  
 269 links, 509 parameters), Insurance (27 nodes, 52 links, 984 parameters), and Barley (48 nodes, 84  
 270 links, 114,005 parameters). We then assessed causal discovery performance in terms of *Precision*,  
 271  $TP/(TP + FP)$ , and *Recall*,  $TP/(TP + FN)$ , relative to the theoretical PAGs, while counting as  
 272 false positive (*FP*), all correctly predicted edges but without or with a different orientation as the  
 273 directed or bidirected edges of the PAG.

274 Fig. 2 compares MIIC\_search&score performance to MIIC results used as starting point for  
 275 MIIC\_search&score and to FCI [41]. MIIC and MIIC\_search&score settings were set as described  
 276 in section 3 above. The open-source MIIC R package (v1.5.2, GPL-3.0 license) was obtained at  
 277 [https://github.com/miicTeam/miic\\_R\\_package](https://github.com/miicTeam/miic_R_package). FCI from the python causal-learn package  
 278 (v0.1.3.8, MIT license) [41] was obtained at <https://github.com/py-why/causal-learn> and  
 279 run with  $G^2$ -conditional independence test and default parameter  $\alpha = 0.05$ .

280 Overall, MIIC\_search&score is found to outperform MIIC in terms of edge precision with little to no  
 281 decrease in edge recall, Fig. 2, demonstrating the benefit of MIIC\_search&score’s rationale to improve  
 282 MIIC predictions by extending MIIC information scores from unshielded triples to higher-order  
 283 information contributions. These originate from *ac*-connected subsets including nodes with more than  
 284 two parents or spouses, or *ac*-connected subsets including two-collider paths. MIIC\_search&score is  
 285 also found to outperform FCI on complex ancestral benchmark networks with many parameters, such  
 286 as Barley (114,005 parameters), Fig. 2. However, FCI is found to reach similar or better precision  
 287 scores on easier benchmarks with fewer parameters (*i.e.* Alarm and Insurance), although its recall  
 288 remains usually lower than MIIC\_search&score, especially at small sample size, as expected for a  
 289 purely constraint-based causal discovery approach.

290 Importantly, the benchmark PAGs used to score the causal discovery results with increasing propor-  
 291 tions of latent variables, Fig. 2, include not only bidirected edges originating from hidden common  
 292 causes but also additional directed or undirected edges arising, in particular, from indirect effects of



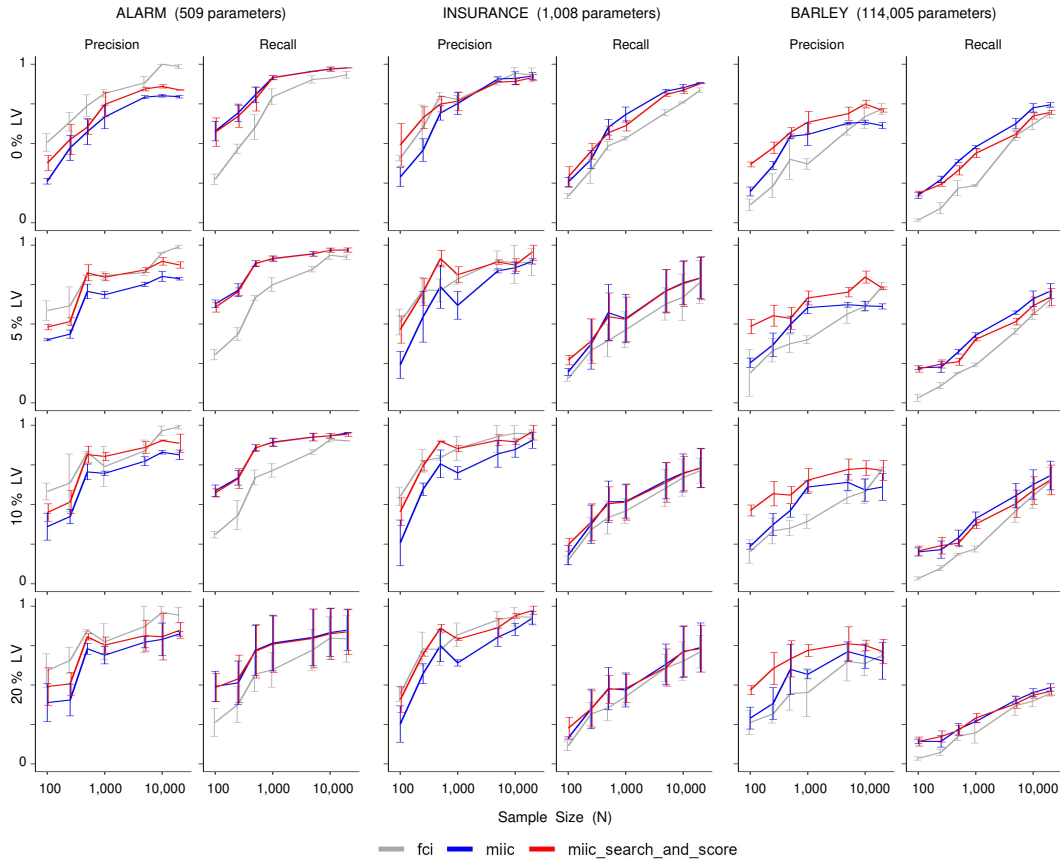


Figure 2: **Benchmark results on ancestral graphs of increasing complexity.** Benchmark results on ancestral graphs obtained by hiding 0%, 5%, 10% or 20% of variables in Bayesian Networks of increasing complexity (see main text): Alarm (lhs), Insurance (middle), and Barley (rhs). MIIC\_search&score results are compared to MIIC results used as starting point for MIIC\_search&score and FCI [41]. Causal discovery performance is assessed in terms of *Precision* and *Recall* relative to the theoretical PAGs, while counting as false positive all correctly predicted edges but without or with a different orientation as the directed or bidirected edges of the PAG. Error bars ( $\pm\sigma$ ): standard deviations.

293 hidden variables with observed parents. Irrespective of their orientations, all these additional edges  
 294 originating from indirect effects of hidden variables generally correspond to weaker effects (*i.e.* lower  
 295 mutual information of indirect effects due to the Data Processing Inequality) and are more difficult to  
 296 uncover than the edges of the original graphical model without hidden variables.

## 297 5 Limitations

298 The main limitation of the paper concerns the local scores used in the search-and-score algorithm,  
 299 which are limited to *ac*-connected subsets of vertices with a maximum of two-collider paths.

300 While this approach could be extended to higher-order information contributions including three-  
 301 or more collider paths, it allows for a simple two-step search-and-score scheme at the level of  
 302 individual nodes (step 1) and edges (step 2), as detailed in section 3. This already shows a significant  
 303 improvement in causal discovery performance (*i.e.* combining good precision and good recall on  
 304 challenging benchmarks) as compared to existing state-of-the-art methods.

## References

- 305 [1] D. Koller, N. Friedman, *Probabilistic Graphical Models: Principles and Techniques* (MIT Press, 2009).
- 306 [2] J. Pearl, A. Paz, Graphoids: A graph-based logic for reasoning about relevance relations, or when would x  
307 tell you more about y if you already know z, *Tech. rep.*, UCLA Computer Science Department (1985).
- 308 [3] J. Pearl, *Probabilistic reasoning in intelligent systems* (Morgan Kaufmann, San Mateo, CA, 1988).
- 309 [4] J. Pearl, *Causality: models, reasoning and inference* (Cambridge University Press, 2009), second edn.
- 310 [5] P. Spirtes, C. Glymour, R. Scheines, *Causation, Prediction, and Search* (MIT press, , 2000), second edn.
- 311 [6] T. S. Richardson, A factorization criterion for acyclic directed mixed graphs, *Proceedings of the Twenty-*  
312 *Fifth Conference on Uncertainty in Artificial Intelligence*, UAI '09 (AUAI Press, Arlington, VA, USA,  
313 2009), p. 462–470.
- 314 [7] J. Tian, J. Pearl, A general identification condition for causal effects, *Proceedings of the National Confer-*  
315 *ence on Artificial Intelligence* (Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999,  
316 2002), pp. 567–573.
- 317 [8] T. Richardson, P. Spirtes, Ancestral graph markov models. *Ann. Statist.* **30**, 962–1030 (2002).
- 318 [9] M. Drton, M. Eichler, T. S. Richardson, Computing maximum likelihood estimates in recursive linear  
319 models with correlated errors. *Journal of Machine Learning Research* **10**, 2329–2348 (2009).
- 320 [10] R. J. Evans, T. S. Richardson, Maximum likelihood fitting of acyclic directed mixed graphs to binary  
321 data. *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence*, UAI'10 (AUAI Press,  
322 Corvallis, OR, USA, 2010).
- 323 [11] S. Triantafillou, I. Tsamardinos, Score-based vs constraint-based causal learning in the presence of  
324 confounders, *CFA@UAI* (2016).
- 325 [12] K. Rantanen, A. Hyttinen, M. Järvisalo, Maximal ancestral graph structure learning via exact search,  
326 *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, C. de Campos,  
327 M. H. Maathuis, eds. (PMLR, 2021), vol. 161 of *Proceedings of Machine Learning Research*, pp. 1237–  
328 1247.
- 329 [13] T. Claassen, I. G. Bucur, Greedy equivalence search in the presence of latent confounders, *Proceedings of*  
330 *the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, J. Cussens, K. Zhang, eds. (PMLR,  
331 2022), vol. 180 of *Proceedings of Machine Learning Research*, pp. 443–452.
- 332 [14] Z. Hu, R. Evans, Faster algorithms for markov equivalence, *Proceedings of the 36th Conference on Uncer-*  
333 *tainty in Artificial Intelligence (UAI)*, J. Peters, D. Sontag, eds. (PMLR, 2020), vol. 124 of *Proceedings of*  
334 *Machine Learning Research*, pp. 739–748.
- 335 [15] W. J. McGill, Multivariate information transmission. *Trans. of the IRE Professional Group on Information*  
336 *Theory (TIT)* **4**, 93-111 (1954).
- 337 [16] H. K. Ting, On the amount of information. *Theory Probab. Appl.* **7**, 439-447 (1962).
- 338 [17] T. S. Han, Multiple mutual informations and multiple interactions in frequency data. *Information and*  
339 *Control* **46**, 26-45 (1980).
- 340 [18] R. W. Yeung, A new outlook on shannon's information measures. *IEEE transactions on information theory*  
341 **37**, 466–474 (1991).
- 342 [19] P. Spirtes, T. Richardson, A polynomial time algorithm for determinint dag equivalence in the presence of  
343 latent variables and selection bias, *Proceedings of the 6th International Workshop on Artificial Intelligence*  
344 *and Statistics* (1996).
- 345 [20] T. Richardson, Markov properties for acyclic directed mixed graphs. *Scandinavian Journal of Statistics* **30**,  
346 145-157 (2003).
- 347 [21] R. A. Ali, T. S. Richardson, Markov equivalence classes for maximal ancestral graphs, *Proceedings of the*  
348 *Eighth Conference on Uncertainty in Artificial Intelligence*, UAI'02 (Morgan Kaufmann Publishers  
349 Inc., San Francisco, CA, USA, 2002), pp. 1–9.
- 350 [22] R. A. Ali, T. S. Richardson, P. Spirtes, J. Zhang, Towards characterizing markov equivalence classes for  
351 directed acyclic graphs with latent variables, *Proceedings of the Fifteenth Conference on Uncertainty in*  
352 *Artificial Intelligence*, UAI'05 (Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2005).
- 353

- 354 [23] J. Tian, Generating markov equivalent maximal ancestral graphs by single edge replacement, *Proceedings*  
355 *of the Fifteenth Conference on Uncertainty in Artificial Intelligence*, UAI'05 (Morgan Kaufmann Publishers  
356 Inc., San Francisco, CA, USA, 2005).
- 357 [24] R. A. Ali, T. S. Richardson, P. Spirtes, Markov equivalence for ancestral graphs. *Ann. Statist.* **37**, 2808–2837  
358 (2009).
- 359 [25] T. Richardson, P. Spirtes, Scoring ancestral graph models, *Tech. rep.* (1999). Available as Technical Report  
360 CMU-PHIL 98.
- 361 [26] J. Zhang, A characterization of markov equivalence classes for directed acyclic graphs with latent variables,  
362 *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, UAI'07 (Morgan  
363 Kaufmann Publishers Inc., San Francisco, CA, USA, 2007).
- 364 [27] J. Zhang, On the completeness of orientation rules for causal discovery in the presence of latent confounders  
365 and selection bias. *Artif. Intell.* **172**, 1873–1896 (2008).
- 366 [28] S. Affeldt, H. Isambert, Robust reconstruction of causal graphical models based on conditional 2-point and  
367 3-point information, *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*,  
368 *UAI 2015, July 12-16, 2015, Amsterdam, The Netherlands* (2015), pp. 42–51.
- 369 [29] L. Verny, N. Sella, S. Affeldt, P. P. Singh, H. Isambert, Learning causal networks with latent variables from  
370 multivariate information in genomic data. *PLoS Comput. Biol.* **13**, e1005662 (2017).
- 371 [30] B. Andrews, G. F. Cooper, T. S. Richardson, P. Spirtes, The m-connecting imset and factorization for admg  
372 models, Preprint (2022). Arxiv 2207.08963.
- 373 [31] Z. Hu, R. J. Evans, Towards standard imsets for maximal ancestral graphs. *Bernoulli* **30** (2024).
- 374 [32] Z. Hu, R. Evans, A fast score-based search algorithm for maximal ancestral graphs using entropy, Preprint  
375 (2024). Arxiv 2402.04777.
- 376 [33] V. Cabeli, H. Li, M. da Câmara Ribeiro-Dantas, F. Simon, H. Isambert, Reliable causal discovery based on  
377 mutual information supremum principle for finite datasets, *WHY21, 35rd Conference on Neural Information*  
378 *Processing Systems* (NeurIPS, 2021).
- 379 [34] M. d. C. Ribeiro-Dantas, H. Li, V. Cabeli, L. Dupuis, F. Simon, L. Hettal, A.-S. Hamy, H. Isambert,  
380 Learning interpretable causal networks from very large datasets, application to 400, 000 medical records of  
381 breast cancer patients. *iScience* **27**, 109736 (2024).
- 382 [35] V. Cabeli, L. Verny, N. Sella, G. Uguzzoni, M. Verny, H. Isambert, Learning clinical networks from medical  
383 records based on information estimates in mixed-type data. *PLoS Comput. Biol.* **16**, e1007866 (2020).
- 384 [36] S. Affeldt, L. Verny, H. Isambert, 3off2: A network reconstruction algorithm based on 2-point and 3-point  
385 information statistics. *BMC Bioinformatics* **17** (2016).
- 386 [37] H. Li, V. Cabeli, N. Sella, H. Isambert, Constraint-based causal structure learning with consistent separating  
387 sets. *Advances in Neural Information Processing Systems (NeurIPS)* **32** (2019).
- 388 [38] P. Kontkanen, P. Myllymäki, A linear-time algorithm for computing the multinomial stochastic complexity.  
389 *Inf. Process. Lett.* **103**, 227–233 (2007).
- 390 [39] T. Roos, T. Silander, P. Kontkanen, P. Myllymäki, Bayesian network structure learning using factorized  
391 nml universal models, *Proc. 2008 Information Theory and Applications Workshop (ITA-2008)* (IEEE Press,  
392 2008).
- 393 [40] M. Scutari, Learning Bayesian Networks with the bnlearn R Package. *J. Stat. Softw.* **35**, 1–22 (2010).
- 394 [41] Y. Zheng, B. Huang, W. Chen, J. Ramsey, M. Gong, R. Cai, S. Shimizu, P. Spirtes, K. Zhang, Causal-learn:  
395 Causal discovery in python. *Journal of Machine Learning Research* **25**, 1–8 (2024).
- 396 [42] Y. M. Shtarkov, Universal sequential coding of single messages. *Problems of Information Transmission* **23**,  
397 3–17 (1987).
- 398 [43] J. Rissanen, I. Tabus, *Adv. Min. Descrip. Length Theory Appl.* (MIT Press, 2005), pp. 245–264.
- 399 [44] W. Szpankowski, *Average case analysis of algorithms on sequences* (John Wiley & Sons, , 2001).

- 400 [45] P. Kontkanen, W. Buntine, P. Myllymäki, J. Rissanen, H. Tirri, Efficient computation of stochastic  
401 complexity. in: *C. Bishop, B. Frey (Eds.) Proceedings of the Ninth International Conference on Artificial  
402 Intelligence and Statistics, Society for Artificial Intelligence and Statistics* **103**, 233–238 (2003).
- 403 [46] P. Kontkanen, Computationally efficient methods for mdl-optimal density estimation and data clustering,  
404 Ph.D. thesis (2009).
- 405 [47] D. M. Chickering, A Transformational Characterization of Equivalent Bayesian Network Structures, *UAI  
406 '95: Proceedings of the Eleventh Annual Conference on Uncertainty in Artificial Intelligence* (Morgan  
407 Kaufmann, 1995), pp. 87–98.
- 408 [48] M. Kalisch, M. Mächler, D. Colombo, M. H. Maathuis, P. Bühlmann, Causal inference using graphical  
409 models with the r package pcalg. *J. Stat. Softw.* **47**, 1–26 (2012).

## 410 Appendix / supplemental material

### 411 A Preliminaries: connection and separation criteria

#### 412 A.1 $m$ -connection vs $m'$ -connection criteria

413 An ancestral graph can be interpreted as encoding a set of conditional independence relations by  
414 a graphical criterion, called  $m$ -separation, based on the concept of  $m$ -connecting paths, which  
415 generalizes the separation criteria of Markov and Bayesian networks to ancestral graphs.

416 **Definition 4.** [ $m$ -connecting path] A path  $\pi$  between  $X$  and  $Y$  is  $m$ -connecting given a (possibly  
417 empty) subset  $C \subseteq V$  (with  $X, Y \notin C$ ) if:

- 418 *i*) its non-collider(s) are not in  $C$ , and
- 419 *ii*) its collider(s) are in  $\text{An}_{\mathcal{G}}(C)$ .

420 **Definition 5.** [ $m$ -separation criterion] The subsets  $A$  and  $B$  are said to be  $m$ -separated by  $C$ , noted  
421  $A \perp_m B | C$ , if there is no  $m$ -connecting path between any vertex in  $A$  and any vertex in  $B$  given  $C$ .

422 The probabilistic interpretation of ancestral graph is given by its global and pairwise Markov properties  
423 (which are equivalent [8]): if  $A$  and  $B$  are  $m$ -separated by  $C$ , then  $A$  and  $B$  are conditionally  
424 independent given  $C$  and  $\forall X \in A$  and  $\forall Y \in B$ , there is a probability distribution  $P$  faithful  
425 to  $\mathcal{G}$  such that their conditional mutual information vanishes, *i.e.*  $I_P(X; Y | C) = 0$ , also noted  
426  $X \perp\!\!\!\perp_P Y | C$ .

427 However, as discussed above, the proof of Theorem 1 will require to introduce a weaker  $m'$ -connection  
428 criterion defined below.

429 **Definition 6.** [ $m'$ -connecting path] A path  $\pi$  between  $X$  and  $Y$  is  $m'$ -connecting given a subset  
430  $C \subseteq V$  (with  $X, Y$  possibly in  $C$ ) if:

- 431 *i*) its non-collider(s) are not in  $C$ , and
- 432 *ii*) its collider(s) are in  $\text{An}_{\mathcal{G}}(\{X, Y\} \cup C)$ .

433 Note, in particular, that an  $m$ -connecting path is necessary an  $m'$ -connecting path but that the  
434 converse is not always true. For example, the path  $X \rightarrow Z \leftarrow T \leftarrow Y$  in Fig. 1G (with  $Z \rightarrow Y$ ) is  
435 an  $m'$ -connecting path given  $T$  (as  $Z \in \text{An}_{\mathcal{G}}(\{X, Y\} \cup T)$ ) but not an  $m$ -connecting path given  $T$   
436 (as  $Z \notin \text{An}_{\mathcal{G}}(T)$ ).

437 However, Richardson and Spirtes 2002 [8] have shown the following lemma,

438 **Lemma 4.** [Corollary 3.15 in [8]] *In an ancestral graph  $\mathcal{G}$ , there is a  $m'$ -connecting path  $\mu$  between*  
439  *$X$  and  $Y$  given  $C$  if and only if there is a (possibly different)  $m$ -connecting path  $\pi$  between  $X$  and*  
440  *$Y$  given  $C$ .*

441 Hence, Lemma 4 implies that  $m'$ -separation and  $m$ -separation criteria are in fact equivalent, as an  
442 absence of  $m'$ -connecting paths implies an absence of  $m$ -connecting paths and vice versa. This  
443 enables to reformulate the  $m$ -separation criterion above as,

444 **Definition 7.** [ $m'$ -separation (and  $m$ -separation) criteria] The subsets  $A$  and  $B$  are said to be  $m'$ -  
445 separated (or  $m$ -separated) by  $C$ , if all paths from any  $X \in A$  to any  $Y \in B$  have either

- 446 *i*) a non-collider in  $C$ , or
- 447 *ii*) a collider *not* in  $\text{An}_{\mathcal{G}}(\{X, Y\} \cup C)$ .

448 The probabilistic interpretation of an ancestral graph is given by its (global) Markov property: if  $A$   
449 and  $B$  are  $m$ -separated (or  $m'$ -separated) by  $C$ , then  $A$  and  $B$  are conditionally independent given  
450  $C$ , noted as,  $A \perp_m B | C$ .

#### 451 A.2 $ac$ -connecting paths and $ac$ -connected subsets

452 Let us now recall the definition of **ancestor collider connecting paths** or  **$ac$ -connecting paths**,  
453 which is directly relevant to characterize the likelihood decomposition and Markov equivalent classes  
454 of ancestral graphs (Theorem 1). We give here a different yet equivalent definition of  $ac$ -connecting  
455 paths as defined in the main text (Definition 2) in order to underline the similarities and differences  
456 with the notion of  $m'$ -connecting path (Definition 6).

457 **Definition 8.** [*ac*-connecting path] A path  $\pi$  between  $X$  and  $Y$  is an *ac*-connecting path given a  
 458 subset  $C \subseteq V$  (with  $X$  and  $Y$  possibly in  $C$ ) if:

- 459 i)  $\pi$  does not have any noncollider, and  
 460 ii) its collider(s) are in  $\text{An}_{\mathcal{G}}(\{X, Y\} \cup C)$ .

461 Hence, more simply (following Definition 2 in the main text), an *ac*-connecting path given  $C$  is a  
 462 collider path,  $X \ast \rightarrow Z_1 \leftrightarrow \dots \leftrightarrow Z_K \leftarrow \ast Y$ , with all  $Z_i \in \text{An}_{\mathcal{G}}(\{X, Y\} \cup C)$ , i.e. with  $Z_i$  in  $C$  or  
 463 connected to  $\{X, Y\} \cup C$  by an ancestor path,  $Z_i \rightarrow \dots \rightarrow T$  with  $T \in \{X, Y\} \cup C$ .

464 **Definition 9.** [*ac*-separation criterion] The subsets  $A$  and  $B$  are said to be *ac*-separated by  $C$  if there  
 465 is no *ac*-connecting path between any vertex in  $A$  and any vertex in  $B$  given  $C$ .

466 Previous definitions and Lemma 4 readily lead to the following corollary between the different  
 467 connection and separation criteria:

468 **Corollary 5.**

- 469 i) *m*-connecting path  $\pi \implies m'$ -connecting path  $\pi$   
 470 ii) *ac*-connecting path  $\pi \implies m'$ -connecting path  $\pi$   
 471 iii) *m*-separation  $\iff m'$ -separation  
 472 iv) *m*/*m'*-separation  $\implies ac$ -separation

473 Finally, we recall the notion of ***ac*-connected subset** (Definition 3 in the main text), which is central  
 474 for the decomposition of the likelihood of ancestral graphs (Theorem 1): A subset  $C$  is said to be  
 475 *ac*-connected if  $\forall X, Y \in C$ , there is an *ac*-connecting path between  $X$  and  $Y$  w.r.t.  $C$ .

## 476 B Proof of Theorem 1.

477 In order to prove that the likelihood function of an ancestral graph, Eq. 12, contains all and only the  
 478 *ac*-connected subsets of vertices in  $\mathcal{G}$  (Definition 3), we will first show (i) that all non-*ac*-connected  
 479 subsets  $S'$  are included in a cancelling combination of multivariate information terms  $I(X; Y | \mathbf{A}) = 0$ ,  
 480 with  $X, Y \in S'$  and  $S' \subseteq S = \{X, Y\} \cup \mathbf{A}$ . Conversely, we will then show (ii) that cancelling  
 481 combinations of multivariate information terms associated to pairwise conditional independence,  
 482  $I(X; Y | \mathbf{A}) = \sum_{S' \subseteq S}^{X, Y \in S'} (-1)^{|S'|} I(S') = 0$  do not contain any *ac*-connected subset  $S'$ . Finally, we  
 483 will prove (iii) that the information terms which appear in multiple cancelling combinations from  
 484 different pairwise independence constraints do not modify the multivariate information decomposition  
 485 of the likelihood function of ancestral graphs, Eq. 12, as these shared/overlapping terms in fact all  
 486 cancel through more global Markov independence relationships involving higher order (three or more  
 487 points) vanishing multivariate information terms, such as  $I(X; Y; Z | \mathbf{A}) = 0$ .

488 i) Let's first prove that all non-*ac*-connected subsets  $S'$  are included in at least one cancelling  
 489 combination of multivariate information terms,  $I(X; Y | \mathbf{A}) = 0$ , with  $X, Y \in S'$  and  $S' \subseteq \{X, Y\} \cup \mathbf{A}$ .

490 If  $S'$  is a non-*ac*-connected subset, there is at least one disconnected pair  $X$  and  $Y$  for which each  
 491 path  $\pi_j$  between  $X$  and  $Y$  contains either some collider(s) not in  $\text{An}_{\mathcal{G}}(S')$  or, if all colliders along  
 492  $\pi_j$  are in  $\text{An}_{\mathcal{G}}(S')$ , there must be some non-collider(s) at node(s)  $Z_j$  but not necessarily in  $S'$ . Let's  
 493 define  $S = S' \cup_j Z_j$ .  $X$  and  $Y$  can be shown to be *m*-separated given  $S \setminus \{X, Y\}$ , as for each  
 494 path  $\pi_j$  between  $X$  and  $Y$ , its non-collider(s) are in  $S$  at node(s)  $Z_j$  (when all collider(s) along  $\pi_j$   
 495 are in  $S'$ ) or there is some collider(s) not in  $\text{An}_{\mathcal{G}}(S')$ , which are not in  $\text{An}_{\mathcal{G}}(S')$  either. The latter  
 496 statement is proven by contradiction assuming that there is a collider at  $Z \notin \text{An}_{\mathcal{G}}(S')$  such that  
 497  $Z \in \text{An}_{\mathcal{G}}(S)$ . There is therefore a directed path  $Z \rightarrow \dots \rightarrow W$  with  $W \in S$ . Hence,  $W \in S'$  or  
 498 there is a noncollider at  $W \in Z_j$  which is on a path  $\pi_j$  between  $X$  and  $Y$  along which all colliders  
 499 are in  $\text{An}_{\mathcal{G}}(S')$  by construction of  $S$ . This leads by induction to  $Z \rightarrow \dots \rightarrow W \rightarrow \dots \rightarrow T$  where  
 500  $T \in S'$  and thus  $Z \in \text{An}_{\mathcal{G}}(S')$ , which is a contradiction. Hence, all non-*ac*-connected subsets  $S'$   
 501 are included in a cancelling combination of multivariate information terms  $I(X; Y | \mathbf{A}) = 0$ , with  
 502  $X, Y \in S'$  and  $S' \subseteq S = \{X, Y\} \cup \mathbf{A}$ .

503 ii) Conversely, we will now show that cancelling combinations of multivariate information terms  
 504 associated to pairwise conditional independence,  $I(X; Y | \mathbf{A}) = \sum_{S' \subseteq S}^{X, Y \in S'} (-1)^{|S'|} I(S') = 0$ , do  
 505 not contain any *ac*-connected subset  $S'$ .

506 We will prove it by contradiction assuming that there exists a subset  $W \subseteq \mathbf{A}$ , such that  $S' =$   
 507  $\{X, Y\} \cup W$  is *ac*-connected. In particular, there should be an *ac*-connecting path between  $X$  and  $Y$

508 confined to  $\mathbf{An}_G(\mathcal{S}')$  and thus to  $\mathbf{An}_G(\mathcal{S}) \supseteq \mathbf{An}_G(\mathcal{S}')$ , which is an  $m'$ -connecting path between  $X$   
509 and  $Y$  given  $\mathbf{A}$ , contradicting the above hypothesis of  $m'$ -separation given  $\mathbf{A}$ , *i.e.*  $I(X; Y|\mathbf{A}) = 0$ .  
510 The use of  $m'$ -separation, *i.e.* the absence of  $m'$ -connecting paths with colliders in  $\mathbf{An}_G(\mathcal{S})$  rather  
511 than  $m$ -connecting paths with colliders in  $\mathbf{An}_G(\mathbf{A})$ , is necessary here, see Definitions 4 and 6. Hence,  
512 no  $ac$ -connected subset  $\mathcal{S}'$  is included in cancelling combinations of multivariate information terms  
513 associated to pairwise conditional independence,  $I(X; Y|\mathbf{A}) = \sum_{\mathcal{S}' \subseteq \mathcal{S}}^{X, Y \in \mathcal{S}'} (-1)^{|\mathcal{S}'|} I(\mathcal{S}') = 0$ .

514 *iii*) Finally, we will show that the information terms which appear in multiple cancelling combina-  
515 tions from different pairwise independence constraints do not modify the multivariate information  
516 decomposition of the likelihood function of ancestral graphs, Eq. 12, as these shared/overlapping  
517 terms in fact all cancel through more global Markov independence relationships involving higher  
518 order (three or more points) vanishing multivariate information terms, such as  $I(X; Y; Z|\mathbf{A}) = 0$ .

519 This result requires to use an ordering of the nodes,  $X_k \succ X_j \succ X_i$ , that is compatible with the  
520 directed edges of the ancestral graph assumed to have no undirected edges, *i.e.*  $X_j \notin \mathbf{An}(X_i)$  if  
521  $X_j \succ X_i$ . Under this ordering, higher order nodes  $X_k \succ X_i \succ X_j$  can be a priori excluded from all  
522 separating sets  $\mathbf{A}_{ij}$  of pairs of lower order nodes, *i.e.* if  $I(X_i; X_j|\mathbf{A}_{ij}) = 0$  then  $X_k \notin \mathbf{A}_{ij}$ .

523 In particular, the two pairwise conditional independence relations  $I(X_k; X_\ell|\mathbf{A}_{k\ell}) = 0$ , with  $X_\ell \succ$   
524  $X_k$ , and  $I(X_i; X_j|\mathbf{A}_{ij}) = 0$ , with  $X_j \succ X_i$ , do not share any multivariate information terms, if  
525  $X_\ell \neq X_j$ . Indeed, as  $I(X_k; X_\ell|\mathbf{A}_{k\ell})$  contains all information terms including both  $X_k$  and  $X_\ell$  as  
526 well as every subset (possibly empty) of  $\mathbf{A}_{k\ell}$ , none of them includes  $X_j$  if  $X_\ell \succ X_j$ . Therefore  
527  $I(X_k; X_\ell|\mathbf{A}_{k\ell})$  does not contain any information term of  $I(X_i; X_j|\mathbf{A}_{ij})$  which contains both  $X_i$  and  
528  $X_j$  as well as every subset (possibly empty) of  $\mathbf{A}_{ij}$ . This property eliminates all multiple counting of  
529 multivariate informations terms shared if  $X_\ell \neq X_j$ . Note that this result does not hold in general for  
530 ancestral graphs including undirected edges.

531 Hence, the issue of redundant multivariate information terms in the likelihood decomposition, Eq. 12,  
532 is related to the conditional independences of two or more pairs,  $\{X_i, X_r\}, \{X_j, X_r\}, \dots, \{X_\ell, X_r\}$ ,  
533 sharing the same higher order node,  $X_r$ . However, this situation also entails a more global Markov  
534 independence constraint between  $X_r$  and  $\{X_i, X_j, \dots, X_\ell\}$ , given a separating set  $\mathbf{A}$ , which can be  
535 decomposed into more local independence constraints using the chain rule and the decomposition  
536 rules of multivariate information (Eq. 9),

$$\begin{aligned}
0 &= I(\{X_i, X_j, \dots, X_\ell\}; X_r|\mathbf{A}) \\
&= (I(X_i; X_r|\mathbf{A}) + I(X_j; X_r|\mathbf{A}, X_i)) + [I(X_k; X_r|\mathbf{A}, X_i, X_j)] + \dots + I(X_\ell; X_r|\mathbf{A}, \dots) \\
&= (I(X_i; X_r|\mathbf{A}) + I(X_j; X_r|\mathbf{A}) - I(X_i; X_j; X_r|\mathbf{A})) \\
&\quad + [I(X_k; X_r|\mathbf{A}, X_i) - I(X_j; X_k; X_r|\mathbf{A}, X_i)] + \dots + I(X_\ell; X_r|\mathbf{A}, \dots) \\
&= (I(X_i; X_r|\mathbf{A}) + I(X_j; X_r|\mathbf{A}) - I(X_i; X_j; X_r|\mathbf{A})) \\
&\quad + [I(X_k; X_r|\mathbf{A}) - I(X_j; X_k; X_r|\mathbf{A}) - I(X_i; X_k; X_r|\mathbf{A}) + I(X_i; X_j; X_k; X_r|\mathbf{A})] + \dots
\end{aligned}$$

537 where all the conditional multivariate information terms vanish by induction due to the non-  
538 negativity of (conditional) mutual information. In particular, the conditional multivariate in-  
539 formation terms in the last expression, *i.e.* between  $X_r$  and each subset of  $\{X_i, X_j, \dots, X_\ell\}$   
540 given the separating set  $\mathbf{A}$ , all vanish. This result can be readily extended to any subsets  
541  $\{X_r, X_s, \dots, X_z\}$  (conditionally) independent of  $\{X_i, X_j, \dots, X_\ell\}$  given a separating set  $\mathbf{A}$ ,  
542 *i.e.*  $I(\{X_i, X_j, \dots, X_\ell\}; \{X_r, X_s, \dots, X_z\}|\mathbf{A}) = 0$ . Hence, as the final conditional multivariate  
543 cross information terms of the decomposition all vanish while not sharing any subsets of variables,  
544 it proves the absence of redundancy and a global cancellation of non- $ac$ -connected subsets (from  
545 pairwise and higher order conditional independence relations) in the likelihood function of ancestral  
546 graphs without undirected edges, Eq. 12.

547 Hence, only  $ac$ -connected subsets effectively contribute to the cross-entropy of an ancestral graph  
548 with only directed and bidirected edges, Eq. 12.  $\square$

## 549 C Factorization of the probability distribution of ancestral graphs

### 550 C.1 Factorization resulting from Theorem 1 and Proposition 3

551 Before presenting the factorization of the model distribution of ancestral graphs resulting from  
 552 Theorem 1 and Proposition 3, it is instructive to obtain an equivalent factorization for Bayesian  
 553 graphs, assuming a positive empirical distributions,  $p(x_1, \dots, x_m) = \prod_{i=1}^m p(x_i | x_{i-1}, \dots, x_1) > 0$ ,

$$\begin{aligned}
 q(x_1, \dots, x_m) &= \prod_{i=1}^m q(x_i | \mathbf{pa}_{x_i}) = \prod_{i=1}^m p(x_i | \mathbf{pa}_{x_i}) \\
 &= p(x_1, \dots, x_m) \prod_{i=1}^m \frac{p(x_i | \mathbf{pa}_{x_i})}{p(x_i | x_{i-1}, \dots, x_1)} \\
 &= p(x_1, \dots, x_m) \prod_{i=1}^m \frac{p(x_i | \mathbf{pa}_{x_i}) p(\mathbf{x}_{i-1} \setminus \mathbf{pa}_{x_i} | \mathbf{pa}_{x_i})}{p(x_i, \mathbf{x}_{i-1} \setminus \mathbf{pa}_{x_i} | \mathbf{pa}_{x_i})} \quad (15)
 \end{aligned}$$

554 This leads to the following alternative expressions for the cross-entropy  $H(p, q) =$   
 555  $-\sum_{\mathbf{x}} p(\mathbf{x}) \log q(\mathbf{x})$  in terms of multivariate entropy and information, which only depend on the  
 556 empirical joint distribution  $p(\mathbf{x})$ ,

$$\begin{aligned}
 H(p, q) &= \sum_{i=1}^m H(x_i | \mathbf{Pa}_{X_i}) \\
 &= H(X_1, \dots, X_m) + \sum_{i=1}^m I(X_i; \mathbf{X}_{i-1} \setminus \mathbf{Pa}_{X_i} | \mathbf{Pa}_{X_i}) \quad (16)
 \end{aligned}$$

557 where  $\sum_{i=1}^m I(X_i; \mathbf{X}_{i-1} \setminus \mathbf{Pa}_{X_i} | \mathbf{Pa}_{X_i})$  can be decomposed, using the chain rule and Eq. 11, into  
 558 unconditional multivariate information terms, which exactly cancel all the multivariate information  
 559 of the non-*ac*-connected subsets of variables in the multivariate entropy decomposition, Eq. 6.

560 Note, however, that this result obtained for Bayesian networks requires an explicit factorization of the  
 561 global model distribution,  $q(\mathbf{x})$ , in terms of the empirical distribution,  $p(\mathbf{x})$ , which is not known and  
 562 presumably does not exist, in general, for ancestral graphs.

563 Alternatively, assuming that the empirical and model distributions are positive ( $\forall \mathbf{x}, p(\mathbf{x}) > 0,$   
 564  $q(\mathbf{x}) > 0$ ), it is always possible to factorize them into factors associated to each (cross) information  
 565 term in the (cross) entropy decomposition, Eq. 6, as,

$$q(\mathbf{x}) = \prod_{i=1}^m q(x_i) \times \prod_{i < j} \frac{q(x_i, x_j)}{q(x_i) q(x_j)} \times \prod_{i < j < k} \frac{q(x_i, x_j, x_k) q(x_i) q(x_j) q(x_k)}{q(x_i, x_j) q(x_i, x_k) q(x_j, x_k)} \times \dots \quad (17)$$

566 where all the marginal distributions over a subset of variables, e.g.  $q(x_i, x_j, x_k) = \sum_{\ell \neq i, j, k} q(\mathbf{x})$  or  
 567  $p(x_i, x_j, x_k) = \sum_{\ell \neq i, j, k} p(\mathbf{x})$ , cancel two-by-two by construction.

568 This can be illustrated on a simple example of a two-collider path including one bidirected edge,  
 569  $X \rightarrow Z \longleftrightarrow Y \leftarrow W$  (Fig. 1D), valid for  $q(\cdot)$  and  $p(\cdot)$  alike,

$$\begin{aligned}
 q(x, z, y, w) &= q(x) q(z) q(y) q(w) \\
 &\times \frac{q(x, z)}{q(x) q(z)} \frac{q(z, y)}{q(z) q(y)} \frac{q(y, w)}{q(y) q(w)} \frac{q(x, y)}{q(x) q(y)} \frac{q(x, w)}{q(x) q(w)} \frac{q(z, w)}{q(z) q(w)} \\
 &\times \frac{q(x) q(z) q(y) q(x, z, y)}{q(x, z) q(x, y) q(z, y)} \frac{q(z) q(y) q(w) q(z, y, w)}{q(z, y) q(z, w) q(y, w)} \\
 &\times \frac{q(x) q(z) q(w) q(x, z, w)}{q(x, z) q(x, w) q(z, w)} \frac{q(x) q(y) q(w) q(x, y, w)}{q(x, y) q(x, w) q(y, w)} \\
 &\times \frac{q(x, z) q(z, y) q(y, w) q(x, y) q(x, w) q(z, w) q(x, z, y, w)}{q(x, z, y) q(x, z, w) q(x, y, w) q(z, y, w) q(x) q(y) q(z) q(w)} \quad (18)
 \end{aligned}$$

570 where all individual distribution marginals on subsets of variables, e.g.  $q(x), q(x, z), q(x, z, y)$  (or  
 571  $p(x), p(x, z), p(x, z, y)$ ), cancel two-by-two by construction, except  $q(x, z, y, w)$  (or  $p(x, z, y, w)$ ).



572 In addition and *only for the model distribution*  $q(\cdot)$ , all ratios in gray in Eq. 18 also cancel due to  
 573 Markov independence relations across non-*ac*-connected subsets (see proof of Theorem 1). This  
 574 leaves a truncated factorization retaining all and only the *ac*-connected subsets of variables in the  
 575 graph, which we propose to estimate on empirical data by substituting the remaining  $q(\cdot)$  terms by  
 576 their empirical counterparts  $p(\cdot)$ , see Proposition 3.

577 This leads to the following global factorization for  $q(\cdot)$  in terms of  $p(\cdot)$ ,

$$\begin{aligned}
 q(x, z, y, w) &\equiv p(x) p(z) p(y) p(w) \frac{p(x, z)}{p(x) p(z)} \frac{p(z, y)}{p(z) p(y)} \frac{p(y, w)}{p(y) p(w)} \\
 &\times \frac{p(x) p(z) p(y) p(x, z, y)}{p(x, z) p(x, y) p(z, y)} \frac{p(z) p(y) p(w) p(z, y, w)}{p(z, y) p(z, w) p(y, w)} \\
 &\times \frac{p(x, z) p(z, y) p(y, w) p(x, y) p(x, w) p(z, w) p(x, z, y, w)}{p(x, z, y) p(x, z, w) p(x, y, w) p(z, y, w) p(x) p(y) p(z) p(w)} \\
 &= p(x, z, y, w) \frac{p(x) p(y)}{p(x, y)} \frac{p(x) p(w)}{p(x, w)} \frac{p(z) p(w)}{p(z, w)} \\
 &\times \frac{p(x, z) p(x, w) p(z, w)}{p(x) p(z) p(w) p(x, z, w)} \frac{p(x, y) p(x, w) p(y, w)}{p(x) p(y) p(w) p(x, y, w)} \tag{19}
 \end{aligned}$$

578 where the terms in gray have been passed to the lhs of Eq. 18 applied to  $p(\cdot)$ . This ultimately  
 579 leads to the analog of the Bayesian Network factorization in Eq. 15 but for the two-collider path,  
 580  $X \rightarrow Z \longleftrightarrow Y \leftarrow W$  (Fig. 1D),

$$q(x, z, y, w) \equiv p(x, z, y, w) \frac{p(x) p(w)}{p(x, w)} \frac{p(z|x) p(w|x)}{p(z, w|x)} \frac{p(x|w) p(y|w)}{p(x, y|w)} \tag{20}$$

581 where the last three factors “correct” the expression of  $p(x, z, y, w)$  for the three (conditional)  
 582 independences entailed by the underlying graph, that is,  $X \perp W$ ,  $Z \perp W|X$ , and  $X \perp Y|W$ .

## 583 C.2 Relation to the head-and-tail factorizations

584 The head-and-tail factorizations of the model distribution of an acyclic directed mixed graph, intro-  
 585 duced by Richardson 2009 [6], enable the parametrization of the joint probability distribution with  
 586 independent parameters for ancestrally closed subsets of vertices.

587 For instance, the head-and-tail factorizations of the simple two-collider path including one bidirected  
 588 edge,  $X \rightarrow Z \longleftrightarrow Y \leftarrow W$ , introduced above, Fig. 1D, are [6],

$$\begin{aligned}
 q(x, w) &= q(x) q(w) \\
 q(x, z) &= q(z|x) q(x) \\
 q(y, w) &= q(y|w) q(w) \\
 q(x, z, w) &= q(z|x) q(x) q(w) \\
 q(x, y, w) &= q(y|w) q(w) q(x) \\
 q(x, z, y, w) &= q(z, y|x, w) q(x) q(w) \tag{21}
 \end{aligned}$$

589 Importantly, these head-and-tail factorizations imply additional relations such as  $q(y|w) = q(y|x, w)$   
 590 (*i.e.*  $X \perp Y|W$ ) obtained by comparing the last two relations in Eq. 21 after marginalizing  
 591  $q(x, z, y, w)$  over  $z$ . However, such implicit conditional independence relations are *not verified*  
 592 *by the empirical distribution*  $p(\cdot)$  in general and prevent the estimation of the head-and-tail factoriza-  
 593 tions by substituting the rhs  $q(\cdot)$  terms in Eq. 21 with their empirical counterparts  $p(\cdot)$ , as in the case  
 594 of Bayesian networks, Eq. 15.

595 Indeed, while the head-and-tail factorization relations, Eq. 21, obey the local and global Markov  
 596 independence relations entailed by the graphical model, Fig. 1D, leading to the cancellation of all  
 597 factors associated to non-*ac*-connected subsets in gray in Eq. 18, the remaining head-and-tail factors  
 598 cannot be readily estimated with the empirical distribution  $p(\cdot)$ .

599 In particular, the cross-entropy of the two-collider path of interest, Fig. 1D, obtained with the head-  
600 and-tail factorizations corresponds to<sup>1</sup>  $H(p, q) = -\sum p(x, z, y, w) \log q(z, y|x, w) q(x) q(w)$ . Then,  
601 estimating the  $q(\cdot)$  terms with their  $p(\cdot)$  counterparts leads to the cross-entropy of a Bayesian graph,  
602 Fig. 1E, with a different Markov equivalent class than the ancestral graph of interest, Fig. 1D. A  
603 similar discrepancy is obtained with a c-component factorization which leads to the cross-entropy of  
604 the Bayesian graph of Fig. 1E without edge  $X \rightarrow Y$ , corresponding to a different Markov equivalence  
605 class than the previous two graphs, Figs. 1D & E.

606 These examples illustrate the difficulty to exploit the c-component or head-and-tail factorizations to  
607 estimate the likelihood of ancestral graphs including bidirected edge(s).

## 608 D Node and edge scores based on Normalized Maximum Likelihood criteria

609 Search-and-score methods based on likelihood estimates need to properly account for finite sample  
610 size, as cross-entropy minimization leads to ever more complex models, resulting in model overfitting  
611 for finite datasets. While BIC regularization is valid in the asymptotic limit of very large datasets, it  
612 tends to overestimate finite size corrections, leading to lower recall, in general. In order to better take  
613 into account finite sample size, we used instead the (universal) Normalized Maximum Likelihood  
614 (NML) criteria [42, 43, 38, 39], which amounts to normalizing the likelihood function over all  
615 possible datasets with the same number  $N$  of samples.

616 **Node score.** We first used the factorized Normalized Maximum Likelihood (fNML) complexity [38,  
617 39] to define a local score for each node  $X_i$ , which extends the decomposable likelihood of Bayesian  
618 graphs given each node's parents, Eq. 2, to all non-descendant neighbors,  $\mathbf{Pa}'_{X_i}$ ,

$$\mathcal{L}_{\mathcal{D}|G_{X_i}} = e^{-N \cdot \text{Score}_n(X_i)} = \frac{e^{-NH(X_i|\mathbf{Pa}'_{X_i})}}{\sum_{|\mathcal{D}'|=N} e^{-NH(X_i|\mathbf{Pa}'_{X_i})}} \quad (22)$$

$$= e^{-NH(X_i|\mathbf{Pa}'_{X_i}) - \sum_j^{q_i} \log C_{n_j}^{r_i}} \quad (23)$$

$$= e^{N \sum_j^{q_i} \sum_k^{r_i} \frac{n_{jk}}{N} \log \left( \frac{n_{jk}}{n_j} \right) - \sum_j^{q_i} \log C_{n_j}^{r_i}} \quad (24)$$

$$= \prod_j \frac{\prod_k^{r_i} \binom{n_{jk}}{n_j}^{n_{jk}}}{C_{n_j}^{r_i}} \quad (25)$$

619 where  $n_{jk}$  corresponds to the number of data points for which  $X_i$  is in its  $k$ th state and its non-  
620 descendant neighbors in their  $j$ th state, with  $n_j = \sum_k^{r_i} n_{jk}$ . The universal normalization constant  $C_n^r$   
621 is then computed by summing the numerator over all possible partitions of the  $n$  data points into a  
622 maximum of  $r$  subsets,  $\ell_1 + \ell_2 + \dots + \ell_r = n$  with  $\ell_k \geq 0$ ,

$$C_n^r = \sum_{\ell_1 + \ell_2 + \dots + \ell_r = n} \frac{n!}{\ell_1! \ell_2! \dots \ell_r!} \prod_{k=1}^r \binom{\ell_k}{n}^{\ell_k} \quad (26)$$

623 which can in fact be computed in linear-time using the following recursion [38],

$$C_n^r = C_n^{r-1} + \frac{n}{r-2} C_n^{r-2} \quad (27)$$

624 with  $C_n^1 = 1$  for all  $n$  and applying Eq. 30 below for  $r = 2$ . However, for large  $n$  and  $r$ ,  $C_n^r$   
625 computation tends to be numerically unstable, which can be circumvented by implementing the  
626 recursion on parametric complexity ratios  $\mathcal{D}_n^r = C_n^r / C_n^{r-1}$  rather than parametric complexities  
627 themselves [35] as,

$$\mathcal{D}_n^r = 1 + \frac{n}{(r-2)\mathcal{D}_n^{r-1}} \quad (28)$$

$$\log C_n^r = \sum_{k=2}^r \log \mathcal{D}_n^k \quad (29)$$

<sup>1</sup>Indeed, all terms in Eq. 18 actually cancel two-by-two by construction, *whatever their factorization expression*, except for the remaining joint-distribution over all variables,  $q(x, z, y, w) = q(z, y|x, w) q(x) q(w)$ .

628 for  $r \geq 3$ , with  $\mathcal{C}_n^1 = 1$  and  $\mathcal{C}_n^2 = \mathcal{D}_n^2$ , which can be computed directly with the general formula,  
 629 Eq. 26, for  $r = 2$ ,

$$\mathcal{C}_n^2 = \sum_{h=0}^n \binom{n}{h} \left(\frac{h}{n}\right)^h \left(\frac{n-h}{n}\right)^{n-h} \quad (30)$$

630 or its Szpankowski approximation for large  $n$  (needed for  $n > 1000$  in practice) [44–46],

$$\mathcal{C}_n^2 = \sqrt{\frac{n\pi}{2}} \left( 1 + \frac{2}{3} \sqrt{\frac{2}{n\pi}} + \frac{1}{12n} + \mathcal{O}\left(\frac{1}{n^{3/2}}\right) \right) \quad (31)$$

$$\simeq \sqrt{\frac{n\pi}{2}} \exp\left(\sqrt{\frac{8}{9n\pi}} + \frac{3\pi - 16}{36n\pi}\right) \quad (32)$$

631

632 This leads to the following local score for each node  $X_i$ , which is minimized over alternative  
 633 combinations of non-descendant neighbors,  $\mathbf{Pa}'_{X_i} \subseteq \mathbf{Pa}_{X_i} \cup \mathbf{Sp}_{X_i} \cup \mathbf{Ne}_{X_i}$ , in the first step of the  
 634 local search-and-score algorithm (step 1) detailed in the main text,

$$\text{Score}_n(X_i) = H(X_i | \mathbf{Pa}'_{X_i}) + \frac{1}{N} \sum_j^{q_{X_i}} \log \mathcal{C}_{n_j}^{r_{X_i}} \quad (33)$$

635 **Edge scores.** We then defined several edge scores to optimize the orientation of each edge,  $X - Y$ ,  
 636 given its close surrounding vertices.

637 To this end, we first introduced a local score for node pairs which simply sums the node scores, Eq. 33,  
 638 for each node. The resulting pair scores are listed in Table 2 for unconnected node pairs and for pairs  
 639 of nodes connected by a directed edge, where  $\mathbf{Pa}'_{X \setminus Y} = \mathbf{Pa}_X \cup \mathbf{Sp}_X \setminus Y$  and  $\mathbf{Pa}'_{Y \setminus X} = \mathbf{Pa}_Y \cup \mathbf{Sp}_Y \setminus X$   
 640 with their corresponding combinations of levels,  $q_{y \setminus x}$  and  $q_{x \setminus y}$ .

Table 2: Local scores for node pairs

Pair score	Information	fNML Complexity
$X \not\sim Y$	$H(X   \mathbf{Pa}'_{X \setminus Y}) + H(Y   \mathbf{Pa}'_{Y \setminus X})$	$\frac{1}{N} \left( \sum_j^{q_{x \setminus y}} \log \mathcal{C}_{n_j}^{r_x} + \sum_j^{q_{y \setminus x}} \log \mathcal{C}_{n_j}^{r_y} \right)$
$X \rightarrow Y$	$H(X   \mathbf{Pa}'_{X \setminus Y}) + H(Y   \mathbf{Pa}'_{Y \setminus X}, X)$	$\frac{1}{N} \left( \sum_j^{q_{x \setminus y}} \log \mathcal{C}_{n_j}^{r_x} + \sum_j^{q_{y \setminus x}^{r_x}} \log \mathcal{C}_{n_j}^{r_y} \right)$
$X \leftarrow Y$	$H(X   \mathbf{Pa}'_{X \setminus Y}, Y) + H(Y   \mathbf{Pa}'_{Y \setminus X})$	$\frac{1}{N} \left( \sum_j^{q_{x \setminus y}^{r_y}} \log \mathcal{C}_{n_j}^{r_x} + \sum_j^{q_{y \setminus x}} \log \mathcal{C}_{n_j}^{r_y} \right)$

641 Then, edge scores for directed edges,  $X \rightarrow Y$  and  $Y \rightarrow X$ , are defined w.r.t. to the edge removal  
 642 score,  $X \not\sim Y$ , by subtracting the pair scores of unconnected pairs to the pair scores of directed  
 643 edges, leading to the following edge orientation scores,

$$\text{Score}(X \rightarrow Y) = -I(X; Y | \mathbf{Pa}'_{Y \setminus X}) + \frac{1}{N} \left( \sum_j^{q_{y \setminus x}^{r_x}} \log \mathcal{C}_{n_j}^{r_y} - \sum_j^{q_{y \setminus x}} \log \mathcal{C}_{n_j}^{r_y} \right) \quad (34)$$

$$\text{Score}(Y \rightarrow X) = -I(X; Y | \mathbf{Pa}'_{X \setminus Y}) + \frac{1}{N} \left( \sum_j^{q_{x \setminus y}^{r_y}} \log \mathcal{C}_{n_j}^{r_x} - \sum_j^{q_{x \setminus y}} \log \mathcal{C}_{n_j}^{r_x} \right) \quad (35)$$

644 However, if  $r_x \neq r_y$ , the fNML complexities of these orientation scores are not identical for  
 645 Markov equivalent edge orientations between nodes sharing the same parents (or spouses) [47],  
 646  $\mathbf{Pa}'_{Y \setminus X} = \mathbf{Pa}'_{X \setminus Y} = \mathbf{Pa}'$  and  $q_{y \setminus x} = q_{x \setminus y}$ , despite sharing the same conditional mutual information,

$$I(X; Y | \mathbf{Pa}') = \frac{1}{2} \left( H(X | \mathbf{Pa}') + H(Y | \mathbf{Pa}', X) \right) + \frac{1}{2} \left( H(X | \mathbf{Pa}', Y) + H(Y | \mathbf{Pa}') \right) \quad (36)$$

647 This suggests to symmetrize the fNML complexities for edge orientation scores by averaging them  
 648 over each directed orientation, as for the conditional information in Eq. 36, leading to the proposed  
 649 fNML complexity for directed edges given in Table 1 in the main text.

650 For bidirected edges, the proposed local orientation score accounts for all  $ac$ -connected subsets in  
651 close vicinity of the bidirected edge, which concerns all subsets including either  $X$  and any combi-  
652 nation (possibly void) of parents or spouses different from  $Y$  (*i.e.* corresponding to the information  
653 contributions  $H(X|\mathbf{Pa}'_{X,Y})$ ) or  $Y$  and any combination of parents or spouses different from  $X$   
654 (*i.e.* corresponding to the information contributions  $H(Y|\mathbf{Pa}'_{Y,X})$ ) or, else, including both nodes  $X$   
655 and  $Y$  plus any combination of their parents or spouses, corresponding to the following information  
656 contribution,  $-I(X; Y|\mathbf{Pa}'_{X,Y})$ , where  $\mathbf{Pa}'_{X,Y} = \mathbf{Pa}'_{X,Y} \cup \mathbf{Pa}'_{Y,X}$ . This last term,  $-I(X; Y|\mathbf{Pa}'_{X,Y})$ ,  
657 contains all the remaining information contributions once the bidirected orientation score is given  
658 relative to the edge removal score (Table 2) as for the two directed orientation scores, above. Finally,  
659 the symmetrized fNML complexity associated with a bidirected edge should be computed with  
660 the whole set of conditioning parents or spouses,  $\mathbf{Pa}'_{X,Y}$ , as indicated in Table 1. Note that this  
661 bidirected orientation score becomes also Markov equivalent to the two directed orientation scores,  
662 as required, when the nodes share the same parents and spouses, *i.e.*  $\mathbf{Pa}'_{X,Y} = \mathbf{Pa}'_{Y,X} = \mathbf{Pa}'_{X,Y}$  and  
663  $q_{x,y} = q_{y,x} = q_{x,y}$  in Table 1.

## 664 E Toy models

665 Fig. 3 shows three simple ancestral models used to test MIIC\_search&score orientation scores  
666 (Table 1) to effectively predict bidirected orientations when the end nodes do not share the same  
667 parents (Model 1), share some parents (Model 2) or when the bidirected edge is part of a longer than  
668 two-collider paths (Model 3).

669 The data is generated from the theoretical DAG using the `rmvDAG` function in the `pcalg` package  
670 [48]. Each node follows a normal distribution, and the data is discretized using `bnlearn`'s discretize  
671 function using Hartemink's pairwise mutual information method [40]. For these toy models, the edge  
672 orientation scores are computed assuming the correct parents of each node.

673 The prediction of the edge orientation scores are summarized in Table 3 in % of replicates displaying  
674 directed edges (wrong) or bidirected edge (correct) as a function of increasing dataset size  $N$ .

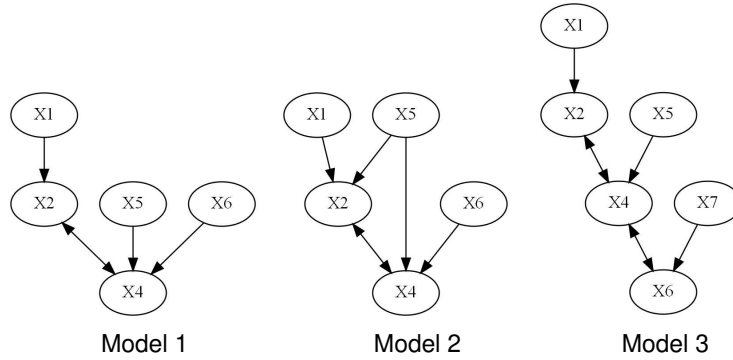


Figure 3: Simple ancestral graphs.

$N$	$\leftarrow$	$\rightarrow$	$\leftrightarrow$	$\leftarrow$	$\rightarrow$	$\leftrightarrow$	$\leftarrow$	$\rightarrow$	$\leftrightarrow$	$\leftarrow$	$\rightarrow$	$\leftrightarrow$
1000	0	100	0	50	42	8	8	88	4	91.7	6.2	2.1
5000	0	68	32	18	2	80	2	80	18	76	24	0
10000	0	10	90	0	0	100	0	6	94	62	22	16
20000	0	0	100	0	0	100	0	0	100	2	0	98
35000	0	0	100	0	0	100	0	0	100	0	0	100
50000	0	0	100	0	0	100	0	0	100	0	0	100

## 675 **NeurIPS Paper Checklist**

### 676 **1. Claims**

677 Question: Do the main claims made in the abstract and introduction accurately reflect the  
678 paper's contributions and scope?

679 Answer: [Yes]

680 Justification: The main claims of the paper are supported by the theoretical and experimental  
681 results shown in Figs. 1 & 2, respectively.

682 Guidelines:

- 683 • The answer NA means that the abstract and introduction do not include the claims  
684 made in the paper.
- 685 • The abstract and/or introduction should clearly state the claims made, including the  
686 contributions made in the paper and important assumptions and limitations. A No or  
687 NA answer to this question will not be perceived well by the reviewers.
- 688 • The claims made should match theoretical and experimental results, and reflect how  
689 much the results can be expected to generalize to other settings.
- 690 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
691 are not attained by the paper.

### 692 **2. Limitations**

693 Question: Does the paper discuss the limitations of the work performed by the authors?

694 Answer: [Yes]

695 Justification: We have added a Discussion & Limitation section at the end of the paper. The  
696 main limitation of the experimental results is the fact that we did not have sufficient time  
697 to perform many dataset replicates of the benchmark ancestral graphs. While the obtained  
698 statistics already support our main experimental results, we intend to perform more dataset  
699 replicates for the final version of the paper.

700 Guidelines:

- 701 • The answer NA means that the paper has no limitation while the answer No means that  
702 the paper has limitations, but those are not discussed in the paper.
- 703 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 704 • The paper should point out any strong assumptions and how robust the results are to  
705 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
706 model well-specification, asymptotic approximations only holding locally). The authors  
707 should reflect on how these assumptions might be violated in practice and what the  
708 implications would be.
- 709 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
710 only tested on a few datasets or with a few runs. In general, empirical results often  
711 depend on implicit assumptions, which should be articulated.
- 712 • The authors should reflect on the factors that influence the performance of the approach.  
713 For example, a facial recognition algorithm may perform poorly when image resolution  
714 is low or images are taken in low lighting. Or a speech-to-text system might not be  
715 used reliably to provide closed captions for online lectures because it fails to handle  
716 technical jargon.
- 717 • The authors should discuss the computational efficiency of the proposed algorithms  
718 and how they scale with dataset size.
- 719 • If applicable, the authors should discuss possible limitations of their approach to  
720 address problems of privacy and fairness.
- 721 • While the authors might fear that complete honesty about limitations might be used by  
722 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
723 limitations that aren't acknowledged in the paper. The authors should use their best  
724 judgment and recognize that individual actions in favor of transparency play an impor-  
725 tant role in developing norms that preserve the integrity of the community. Reviewers  
726 will be specifically instructed to not penalize honesty concerning limitations.

727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: For the theoretical results (notably Theorem 1) we provide the full set of assumptions (section 2 and Appendix A) and a complete proof (Appendix B).

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We provided the full description of the experiments run in the paper (sections 2 & 3 and Appendix D). The open-source code reproducing the experimental results presented in the paper will be provided with the camera-ready version of the paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

779 (d) We recognize that reproducibility may be tricky in some cases, in which case  
780 authors are welcome to describe the particular way they provide for reproducibility.  
781 In the case of closed-source models, it may be that access to the model is limited in  
782 some way (e.g., to registered users), but it should be possible for other researchers  
783 to have some path to reproducing or verifying the results.

## 784 5. Open access to data and code

785 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
786 tions to faithfully reproduce the main experimental results, as described in supplemental  
787 material?

788 Answer: [No]

789 Justification: We do not include a new code with the initial submission, as it is not yet  
790 properly packaged at submission time, but we definitely intend to release this open-source  
791 code including proper annotation and userguide with the final camera-ready version of the  
792 paper. MIIC and FCI open-source packages used for benchmark comparison are already  
793 published and available on public servers.

794 Guidelines:

- 795 • The answer NA means that paper does not include experiments requiring code.
- 796 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
797 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 798 • While we encourage the release of code and data, we understand that this might not be  
799 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
800 including code, unless this is central to the contribution (e.g., for a new open-source  
801 benchmark).
- 802 • The instructions should contain the exact command and environment needed to run to  
803 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
804 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 805 • The authors should provide instructions on data access and preparation, including how  
806 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 807 • The authors should provide scripts to reproduce all experimental results for the new  
808 proposed method and baselines. If only a subset of experiments are reproducible, they  
809 should state which ones are omitted from the script and why.
- 810 • At submission time, to preserve anonymity, the authors should release anonymized  
811 versions (if applicable).
- 812 • Providing as much information as possible in supplemental material (appended to the  
813 paper) is recommended, but including URLs to data and code is permitted.

## 814 6. Experimental Setting/Details

815 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
816 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
817 results?

818 Answer: [Yes]

819 Justification: We provided the full description of the experiments run in the paper (sections  
820 2 3 and Appendix D).

821 Guidelines:

- 822 • The answer NA means that the paper does not include experiments.
- 823 • The experimental setting should be presented in the core of the paper to a level of detail  
824 that is necessary to appreciate the results and make sense of them.
- 825 • The full details can be provided either with the code, in appendix, or as supplemental  
826 material.

## 827 7. Experiment Statistical Significance

828 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
829 information about the statistical significance of the experiments?

830 Answer: [Yes]

831 Justification: The 1-sigma error bars are plotted in Fig. 2. While these statistics already  
832 support our experimental results, we intend to perform more dataset replicates for the  
833 final version of the paper, which we did not have sufficient time to perform by the time of  
834 submission. This should reduce some error bars, in particular, those for the results displaying  
835 large error bars.

836 Guidelines:

- 837 • The answer NA means that the paper does not include experiments.
- 838 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
839 dence intervals, or statistical significance tests, at least for the experiments that support  
840 the main claims of the paper.
- 841 • The factors of variability that the error bars are capturing should be clearly stated (for  
842 example, train/test split, initialization, random drawing of some parameter, or overall  
843 run with given experimental conditions).
- 844 • The method for calculating the error bars should be explained (closed form formula,  
845 call to a library function, bootstrap, etc.)
- 846 • The assumptions made should be given (e.g., Normally distributed errors).
- 847 • It should be clear whether the error bar is the standard deviation or the standard error  
848 of the mean.
- 849 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
850 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
851 of Normality of errors is not verified.
- 852 • For asymmetric distributions, the authors should be careful not to show in tables or  
853 figures symmetric error bars that would yield results that are out of range (e.g. negative  
854 error rates).
- 855 • If error bars are reported in tables or plots, The authors should explain in the text how  
856 they were calculated and reference the corresponding figures or tables in the text.

## 857 8. Experiments Compute Resources

858 Question: For each experiment, does the paper provide sufficient information on the com-  
859 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
860 the experiments?

861 Answer: [Yes]

862 Justification: The computer resource used for all experiments is a simple laptop with intel i7  
863 processors, 12 cores and 16 threads.

864 Guidelines:

- 865 • The answer NA means that the paper does not include experiments.
- 866 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
867 or cloud provider, including relevant memory and storage.
- 868 • The paper should provide the amount of compute required for each of the individual  
869 experimental runs as well as estimate the total compute.
- 870 • The paper should disclose whether the full research project required more compute  
871 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
872 didn't make it into the paper).

## 873 9. Code Of Ethics

874 Question: Does the research conducted in the paper conform, in every respect, with the  
875 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

876 Answer: [Yes]

877 Justification: The paper does not use or produce sensitive data nor concern potentially  
878 harmful applications.

879 Guidelines:

- 880 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 881 • If the authors answer No, they should explain the special circumstances that require a  
882 deviation from the Code of Ethics.



- 883 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
884 eration due to laws or regulations in their jurisdiction).

885 **10. Broader Impacts**

886 Question: Does the paper discuss both potential positive societal impacts and negative  
887 societal impacts of the work performed?

888 Answer: [Yes]

889 Justification: The paper does not use or produce sensitive data nor concern potentially  
890 harmful applications.

891 Guidelines:

- 892 • The answer NA means that there is no societal impact of the work performed.
- 893 • If the authors answer NA or No, they should explain why their work has no societal  
894 impact or why the paper does not address societal impact.
- 895 • Examples of negative societal impacts include potential malicious or unintended uses  
896 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
897 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
898 groups), privacy considerations, and security considerations.
- 899 • The conference expects that many papers will be foundational research and not tied  
900 to particular applications, let alone deployments. However, if there is a direct path to  
901 any negative applications, the authors should point it out. For example, it is legitimate  
902 to point out that an improvement in the quality of generative models could be used to  
903 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
904 that a generic algorithm for optimizing neural networks could enable people to train  
905 models that generate Deepfakes faster.
- 906 • The authors should consider possible harms that could arise when the technology is  
907 being used as intended and functioning correctly, harms that could arise when the  
908 technology is being used as intended but gives incorrect results, and harms following  
909 from (intentional or unintentional) misuse of the technology.
- 910 • If there are negative societal impacts, the authors could also discuss possible mitigation  
911 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
912 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
913 feedback over time, improving the efficiency and accessibility of ML).

914 **11. Safeguards**

915 Question: Does the paper describe safeguards that have been put in place for responsible  
916 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
917 image generators, or scraped datasets)?

918 Answer: [NA]

919 Justification: The paper does not use or produce sensitive data nor concern potentially  
920 harmful applications.

921 Guidelines:

- 922 • The answer NA means that the paper poses no such risks.
- 923 • Released models that have a high risk for misuse or dual-use should be released with  
924 necessary safeguards to allow for controlled use of the model, for example by requiring  
925 that users adhere to usage guidelines or restrictions to access the model or implementing  
926 safety filters.
- 927 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
928 should describe how they avoided releasing unsafe images.
- 929 • We recognize that providing effective safeguards is challenging, and many papers do  
930 not require this, but we encourage authors to take this into account and make a best  
931 faith effort.

932 **12. Licenses for existing assets**

933 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
934 the paper, properly credited and are the license and terms of use explicitly mentioned and  
935 properly respected?

936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986

Answer: [Yes]

Justification: We have credited all previously published resources (including license details) used in the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We do not include a new code with the initial submission, as it is not yet properly packaged at submission time, but we definitely intend to release this open-source code including proper annotation and userguide with the final camera-ready version of the paper.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005

**15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.