# Tools and Algorithms for Causally Interpreting Directed Edges in Maximal Ancestral Graphs

Giorgos Borboudakis, Sofia Triantafillou, Ioannis Tsamardinos
Computer Science Department, University of Crete, Greece
Institute of Computer Science, FORTH, Greece
borbudak@ics.forth.gr

## Abstract

The Maximal Ancestral Graph (MAG) formalism is an important generalization of Bayesian Networks for representing causal processes that admit the possibility of latent confounding variables. Thus, when learning MAGs from data for Causal Discovery, the often unrealistic assumption of Causal Sufficiency can be dismissed. However, the *causal interpretation* of edges in a MAG is not trivial and it is potentially misleading to unfamiliar practitioners. An edge $X \to Y$ may denote either (a) $X$ causes $Y$ and no latent confounding variable is present (*pure-causal* edge) or (b) $X$ causes $Y$ with the potential presence of a latent common cause. In addition, an edge $X \to Y$ may denote (I) $X$ causes $Y$ directly (*direct-causal* edge), i.e., without any modeled variables mediating the causation or (II) $X$ causes $Y$ possibly-indirectly. In this paper, we present polynomial-time algorithms and tools that can distinguish among the above cases and facilitate the causal interpretation of MAGs. In addition, we run simulated experiments to quantify the percentage of edges that can be labeled as pure or direct-causal. Our results show that the percentage of edges that can be labeled as pure-causal achieves a minimum for sparse or dense networks, and a maximum for in-between values of edge density. In contrast, the percentage of edges that can be labeled as direct-causal decreases as the edge density of the MAG increases.

## 1 Introduction

A Causal Bayesian Network (CBN) (Pearl, 2000) is a probabilistic graphical model that can represent a data distribution. The graph of a CBN encodes a set of conditional independencies on the observed variables provided by the *d*-separation criterion. In addition, an edge $X \to Y$ in the graph has *causal* semantics: [$X$ *directly causes* $Y$], that is an intervention of $X$ has a direct effect on $Y$. The meaning of "direct" causation is that the causation persists even when all other *observed* variables are held constant, for some appropriate combination of values of the other variables (Spirtes et al., 2000). Thus, direct causation depends on the context of observed variables.

A major assumption of CBNs is that there are no latent confounding variables, i.e,. latent common causes of the observed variables,

named the *Causal Sufficiency assumption*. If this is not the case, then a system may not be representable so that both the probabilistic *and* the causal semantics are correct. For example, let us assume that the true causal structure is $X \leftarrow H \to Y$, $H$ is latent, and $X$ and $Y$ are dependent. There are three possible CBNs. The graph [$X$ (no edge) $Y$], entails $X$ is independent of $Y$ and does not comply with the probabilistic semantics. The graphs [$X \leftarrow Y$] and [$X \to Y$] correctly represent the distribution but do not comply with the causal semantics.

The Maximal Ancestral Graph (MAG) formalism (Richardson and Spirtes, 2002) is an extension of CBNs which admits the presence of hidden confounders, without explicitly introducing them in the model. An additional important feature of MAGs is that they are closed under marginalization; given a MAG $M$, the

MAG $M'$ representing the causal relations after marginalizing out a set of variables can be computed. The probabilistic interpretation of the edges also stems from a simple generalization of the $d$-separation criterion: the $m$-separation. The semantics of an edge $X \leftrightarrow Y$ are that [neither $X$ causes $Y$ nor the reverse]. An edge $X \leftrightarrow Y$ may be required in the model to represent possible latent confounders: $X \leftarrow H \rightarrow Y$ can be encoded as the MAG $X \leftrightarrow Y$.

However, the advantages of MAGs come at the price of *difficulty in the causal interpretation of the directed edges*. The causal semantics of a directed edge $X \rightarrow Y$ are that [$X$ causes $Y$]. The causal relation may be **direct-causal** (no observed variables mediate it) or it may be **possibly-indirect** and mediated by the *observed* variables. In addition, the edge $X \rightarrow Y$ does not distinguish between two possible scenarios: (a) $X$ causes $Y$ but there also is a latent confounder $H$ of the two variables (thus, the graph is [$X \rightarrow Y \leftarrow H \rightarrow X$]), called a **possibly-confounded** edge, and (b) $X$ causes $Y$ and there are no latent confounders. We call the later a **pure-causal** edge (a *visible edge* in the terminology of (Zhang, 2008)).

In this paper, based on results from (Zhang, 2008), we present algorithms that identify all the direct-causal and pure-causal directed edges in a MAG. The algorithms are complete in the sense that if an edge is labeled as possibly-indirect then it could be the marginal of a DAG where it is mediated by observed variables; similarly, if it is labeled as possibly-confounded it could be the marginal of a DAG containing a marginalized common cause of the edge-points.

We consider labeling **pure-causal** and **direct-causal** edges important for interpreting the graph. For a confounded edge $X \rightarrow Y$ the correlation (equivalently dependency, mutual information) between $X$ and $Y$ is partly due both to the direct causal effect as well as due to the latent confounder (as well as other paths). When $X \rightarrow Y$ is pure-causal edges there are no latent confounders to contribute to the correlation, which makes identification of the total or direct effect, easier (see (Cai et al., 2008) for more details). Similarly, a direct-causal edge

$X \rightarrow Y$ implies that the causal effect of $X$ persists even when all other observed variables are held fixed (for some appropriate set of their values).

Figure 1 shows an example encompassing all concepts and cases of this paper. Figure 1(c) shows a MAG annotated by our tools. Directed edges are partitioned into four categories for the four combinations of pure or possibly-confounded causal edge, and direct or possibly-indirect causal edge. Figure 1(a) and Figure 1(b) show two DAGs that could produce the MAG after marginalizing the $H$ variables. Finally, in simulated experiments we show that the percentage of edges that can be labeled as pure-causal or direct-causal depends on the density of the graph.

## 2 Background

We briefly review the Maximal Ancestral Graph (MAG) formalism (Richardson and Spirtes, 2002) and basic background concepts. MAGs can deal with selection variables but this is out of the scope of this paper (i.e., we actually consider the sub-class called *Directed Maximal Ancestral Graphs* (DMAGs)).

We denote a variable with an upper-case letter (e.g. $X$) and a set of variables with upper-case bold letters (e.g. $\mathbf{Z}$). A *directed mixed graph* $G = (V, E)$ consists of a set of vertices $V$ and a set of edges $E$. We denote with $n = |V|$ the number of vertices and with $m = |E|$ the number of edges of $G$. Directed mixed graphs contain two kinds of edges: directed edges ($\rightarrow$) and bi-directed edges ($\leftrightarrow$). Directed edges model causal ancestry relations, whereas bi-directed edges denote that neither variable causes the other. Each edge has two *marks* (or *orientations*), tails (-) and/or arrowheads ($>$). A wildcard mark ($*$) can be a tail or an arrowhead. An edge $X* \rightarrow Y$ is *into* $Y$, whereas an edge $X \rightarrow Y$ is *out of* $X$.

A *path* $p$ in a directed mixed graph $G$ is a sequence of distinct vertices $p = \langle V_1, V_2, \ldots, V_{|p|} \rangle$, s.t. for $1 \leq i < |p|$, $V_i$ and $V_{i+1}$ are adjacent in $G$. The first and last vertices on a path are called the *end-points* of the path. $X$ is a *parent*

(a) DAG $G_1$

(b) DAG $G_2$

(c) DMAG $M$ after marginalizing variable $H$ in $G_1$ or $H_1$ and $H_2$ in $G_2$ using the Marginalization Procedure
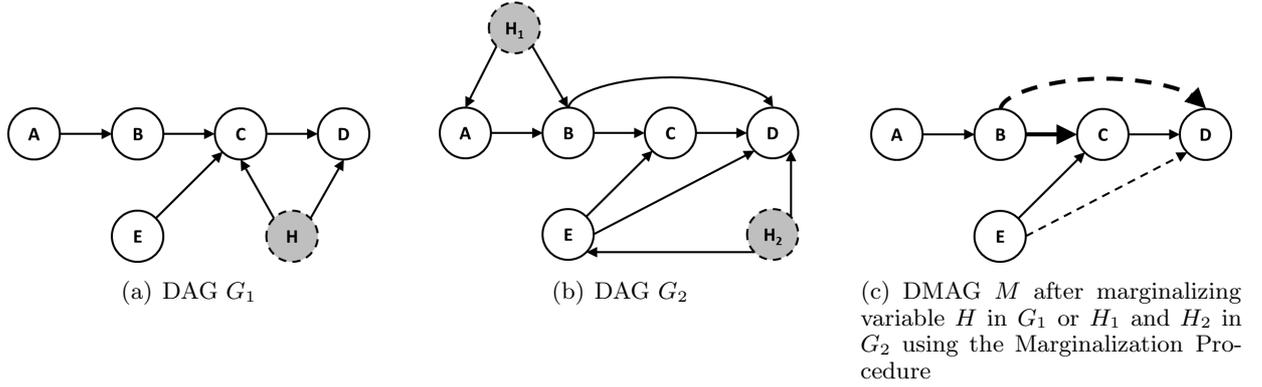
Figure 1: Directed edges in $M$ are annotated: bold edges denote pure-causal relations; dashed edges denote possibly-indirect relations. For example, edge $A \rightarrow B$ in $M$ is possibly-confounded (not bold) since there exist an extension DAG ($G_2$) where $A \leftarrow H \rightarrow B$. Edge $B \rightarrow D$ in $M$ is possibly-indirect (dashed) since the edge is missing from $G_1$. Edge $E \rightarrow D$ in $M$ is both possibly-indirect and possibly-confounded: it is missing from $G_1$ and it is confounded in $G_2$. Edge $B \rightarrow C$ is both a pure-causal and a direct-causal edge: there is no extension DAG such that $B$ and $C$ are confounded or the observed variables mediate the causation.

of $Y$ and $Y$ a *child* of $X$ if $X \rightarrow Y$ is in $G$. A path $p$ is *directed* if for $1 \leq i < |p|$, $V_i$ is a parent of $V_{i+1}$ in $G$. $X$ is an *ancestor* of $Y$ and $Y$ a *descendant* of $X$ if there is a directed path from $X$ to $Y$ in $G$. A *directed cycle* occurs in $G$ if $Y \rightarrow X$ and $X$ is an ancestor of $Y$. An *almost directed cycle* occurs in $G$ if $Y \leftrightarrow X$ and $X$ is an ancestor of $Y$. A triple $\langle X, Y, Z \rangle$ is said to form a *collider* if $X$ and $Z$ are into $Y$. A path $p$ is called a *collider path* if every subsequent triple $\langle V_{i-1}, V_i, V_{i+1} \rangle$, $1 < i < |p|$ in $p$ is a collider.

**Definition 1** (Inducing Path). A path $p$ is *inducing* relative to a set of vertices $\mathbf{L}$ if: (a) every non-endpoint vertex on $p$ is either in $\mathbf{L}$ or a collider, and (b) every collider on $p$ is an ancestor of an end-point vertex of the path. If the set of vertices $\mathbf{L}$ is empty, the path is called a *primitive* inducing path.

**Definition 2** (Directed Maximal Ancestral Graphs). A directed mixed graph is called a Directed Maximal Ancestral Graph (DMAG) if: (a) the graph does not contain any directed cycles, (b) the graph does not contain any almost directed cycles, and (c) for every pair of non-adjacent vertices in the graph, there is no primitive inducing path between them.

First, we want to point out that the definition

of inducing paths presented here differs from the standard definition of inducing paths (Richardson and Spirtes, 2002), since we don't deal with selection variables in this paper. The first two conditions imply that in a DMAG *an arrowhead denotes non-ancestry*. If $X$ and $Y$ are adjacent, and $X$ is into $Y$, $Y$ cannot be an ancestor of $X$ because it would violate one of the conditions (if $X \rightarrow Y$, there would be a directed cycle, if on the other hand $X \leftrightarrow Y$, there would be an almost directed cycle). The third condition implies that two vertices are adjacent if and only if there is a primitive inducing path between them (Richardson and Spirtes, 2002). Notice that all conditions are met by Directed Acyclic Graphs (DAGs). Thus, DAGs are DMAGs without bi-directed edges. A DMAG $M$ over variables $\mathbf{O} \setminus \mathbf{L}$ can be constructed by a DAG (or DMAG) $G$ over variables $\mathbf{O}$ and a set of latent variables $\mathbf{L}$ as follows (Zhang, 2008):

**Procedure 1** (Marginalization). (1) $M$ contains the same variables as $G$ except the latent variables $\mathbf{L}$, (2) two variables $X$ and $Y$ are adjacent in $M$ if and only if there is an inducing path between them relative to $\mathbf{L}$ in $G$, and (3) for two adjacent variables $X$ and $Y$ in $M$, the edge is oriented as (i) $X \rightarrow Y$ if $X$ is an ancestor of $Y$ in $G$, (ii) $X \leftarrow Y$ if $Y$ is an ancestor of

$X$ in $G$ and, (iii) $X \leftrightarrow Y$ otherwise.

Figure 1(c) shows the marginal of the DAGs in Figure 1(a) or 1(b), where $L = \{H\}$ and $L = \{H_1, H_2\}$ respectively. It can be shown that a DMAG $M$ created by this procedure represents the independencies of the marginal distribution entailed by the given $G$ (Richardson and Spirtes, 2002) and additionally represents the causal semantics of $G$ (since it retains the ancestral relations). We call $M$ the **marginal graph** of $G$ over $\mathbf{L}$ and $G$ an **extension** of $M$.

## 3 Problem Definition

Different DAGs may be represented by the same marginal DMAG over some observed variables for different sets of latent variables $\mathbf{L}$. Examples are shown in Figure 1. While the Marginalization Procedure ensures that the resulting DMAG represents the same set of conditional independencies among the observed variables, it does not distinguish between different *causal interpretations* of the edges. To facilitate the causal interpretation we define and subsequently solve the following problems:

**Problem 1.** Given a DMAG $M$, classify all directed edges $X \to Y$ into the following categories: **Pure-causal** if there is no extension DAG $G$ of $M$ such that $X \leftarrow H \to Y$, for some latent variable $H$, and **Possibly-confounded** if there is an extension DAG $G$ where $X \leftarrow H \to Y$ for some latent variable $H$.

We note that (see Section 4 for more details), when an edge is possibly-confounded there is no way of determining whether it is actually confounded or not, solely by examining the graph (there may be other ways to induce the presence of possible confounders when assuming specific functional relations and noise distributions, e.g., (Peters et al., 2012)).

Another problem is that sometimes an edge may be present in a DMAG, while not being present in an extension DAG. Thus, this edge does not represent a direct causal relation in the context of the observed variables. This could happen because of inducing paths: two vertices are adjacent if there is an inducing path between them. We now define the following problem:
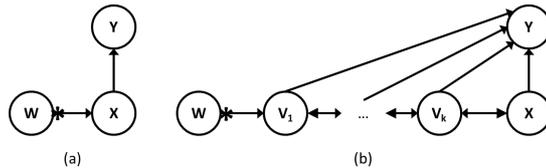


Figure 2: All cases of pure-causal edges. The "*" denotes that the end-point can be either tail or arrowhead. If any of the above induced subgraphs is present in a DMAG, edge $X \to Y$ is pure-causal.

**Problem 2.** Given a DMAG $M$, classify all directed edges $X \to Y$ into the following categories: **Direct-causal** if the edge is present in all extension DAGs of $M$, and **Possibly-indirect** if there is an extension DAG $G$ of $M$ that does not include the edge.

Direct-causal edges $X \to Y$ are important because they imply that the causal effect of $X$ is present even when all other observed variables are held fixed (for some appropriate joint combination of their values). When an edge is possibly-indirect it is impossible to determine, solely by analyzing the graph, whether it is actually indirect or not.

## 4 Labeling Pure-causal and Possibly-Confounded Edges

The theory to distinguish between pure-causal and possibly-confounded directed edges is given in (Zhang, 2008). Actually, in (Zhang, 2008) the concepts of *visible* and *invisible* edges are defined, and subsequently, it is proved that they correspond to our definitions of pure-causal and possibly-confounded edges. The definition now follows:

**Definition 3** (Visibility). A directed edge $X \to Y$ is *visible* if there is a vertex $W$ not adjacent to $Y$, such that either there is an edge between $W$ and $X$ that is into $X$, or there is a collider path between $W$ and $X$ that is into $X$ and every vertex on that path is a parent of $Y$. Otherwise $X \to Y$ is said to be *invisible*.

All possible cases of visible (pure-causal) edges are shown in Figure 2. (Zhang, 2008)

**Algorithm 1** Find pure-causal Edges(DMAG $\mathcal{M}$)

---

1: Mark all directed edges as possibly-confounded
2: $Q \leftarrow$ Queue()
3: **for** each $X \rightarrow Y \in \mathcal{M}$ **do**
4:    **if** $\exists W$ s.t. $W* \rightarrow X \wedge \neg adjacent(W, Y)$ **then**
5:       Mark $X \rightarrow Y$ as pure-causal
6:       $Q$.enqueue($\{X, Y\}$)
7:    **end if**
8: **end for**
9: **while** $\neg Q$.isEmpty() **do**
10:    $\{X, Y\} = Q$.dequeue()
11:    **for** each $V$ s.t. $X \leftrightarrow V \wedge V \rightarrow Y$ **do**
12:       **if** $V \rightarrow Y$ is possibly-confounded **then**
13:          Mark $V \rightarrow Y$ as pure-causal
14:          $Q$.enqueue($\{V, Y\}$)
15:       **end if**
16:    **end for**
17: **end while**
18: **return** All Marked Edges

---

proves that if a directed edge $X \rightarrow Y$ is visible then there is no inducing path between $X$ and $Y$ relative to the set of hidden variables that is into $X$. Consequently, $X$ and $Y$ cannot share a hidden common cause. Thus visibility is sufficient for $X \rightarrow Y$ to be pure-causal. Note that this does not imply that there are no other inducing paths between $X$ and $Y$ and so a visible edge can still be indirect: there can still be inducing paths relative to other sets of variables. Another important fact established by (Zhang, 2008) is that if $X \rightarrow Y$ is invisible, there is at least one DAG extension of the given DMAG, in which $X$ and $Y$ share a hidden common cause. As a result, visibility is also necessary for $X$ and $Y$ to be pure-causal. Thus, the two concepts are equivalent and no further distinction is necessary.

Algorithm 1 finds all pure-causal edges in a DMAG. We will prove that it is sound and complete and provide an upper bound for its time complexity.

**Theorem 1.** *Algorithm 1 is **sound** and **complete**.*

*Proof.* We will call a pure-causal path of length $j$ for the edge $V_j \rightarrow Y$ a path of the form $W* \rightarrow V_1 \leftrightarrow \ldots \leftrightarrow V_j$ and $V_i \rightarrow Y$, for all $i = 1, \ldots, j$. As discussed, an edge $V_j \rightarrow Y$ is pure-causal if and only if there exists a pure-causal path of any length. We will use induction on the maximum length $k$ of a causal path. For $k = 1$, the first For loop of the algorithm identifies all pure-causal edges by enumeration; these are the cases in Figures 2(a) and 2(b). Let us assume (inductive hypothesis) that the algorithm marks as pure-causal all edges $V_j \rightarrow Y$ that have a pure-causal path of length at most $k$. In the While loop, an edge $V_{j+1} \rightarrow Y$ will be marked as pure-causal if and only if (a) the subgraph $V_j \leftrightarrow V_{j+1} \rightarrow Y$ is present, and (b) $V_j \rightarrow Y$ is in the queue and thus has been marked pure-causal. By the inductive hypothesis and (b) above, there exist a pure-causal path of length at most $k$. By (a) this path can be extended to a *sequence* $W* \rightarrow V_1 \leftrightarrow \ldots \leftrightarrow V_j \leftrightarrow V_{j+1}$, s.t., $V_i \rightarrow Y$, for all $i = 1, \ldots j+1$. If the sequence is not a path, then $V_{j+1}$ coincides with some other node $V_i$, $1 \leq i \leq j$, and thus it has a pure-causal path of length strictly less than $k + 1$. In that case, it would have already been marked and would not have pass the test at Line 12. In any case, there exist a causal-path of length at most $k + 1$ for the marked edge $V_{j+1} \rightarrow Y$. $\square$

**Theorem 2.** *Algorithm 1 has $O(n \cdot m)$ time complexity.*

*Proof.* The computational costly parts of the algorithm are at lines 3-8 and 9-17. Note that each operation on the queue takes constant time $O(1)$.

The loop at line 3 will run $O(m)$ times. For the if statement at line 4, we have to find all vertices $W$, s.t. $W \rightarrow X$ or $W \leftrightarrow X$ is in the given DMAG, which can be computed in $O(n)$ time (each vertex is adjacent to at most $n - 1$ other vertices), and then out of those check if any is not adjacent to $Y$, which takes $O(1)$ time for fixed vertex $W$. The statements at lines 5 and 6 take constant time. Thus, lines 3-8 take $O(n \cdot m)$ time.

At this point, $Q$ may contain at most $O(m)$ elements. At each iteration we remove one ele-

ment (line 10). Line 11 takes $O(n)$ time (similar to line 4). Lines 12 and 13 take constant time. The only line that could cause a problem is line 14, because it increases the number of elements in $Q$ and could lead to additional iterations. Notice however that an element is added to the queue only if a pure-causal edge is found which was previously marked as possibly-confounded (line 12). But if it was possibly-confounded, it never was in $Q$. Therefore $Q$ will contain each pure-causal edge **exactly once**. But the number of pure-causal edges is at most $O(m)$. Thus, the outer loop (line 9) will run $O(m)$ times, with a cost of $O(n)$ each, and the total running time of lines 9-17 is $O(n \cdot m)$.

Both, lines 3-8 and 9-17 take $O(n \cdot m)$ time, therefore the algorithm *Find pure-causal Edges* runs in $O(n \cdot m)$ time. ☐

## 5 Definitely Direct and possibly-indirect Edges

An edge $X \to Y$ may not correspond to a direct causal relation, e.g., the edge $B \to D$ in the DMAG shown in Figure 1(c) does not exist in a DAG extension of that DMAG shown in Figure 1(a). Given a DAG $G$ that misses the edge $X \to Y$, a set of latent variables $L$ and the marginal DMAG $M$ relative to $L$, an indirect edge $X \to Y$ appears in $M$ when (a) there is a directed path $X \to \ldots \to Y$ in $G$ and (b) there exists an inducing path relative to $L$ between $X$ and $Y$ in $G$. This follows directly from the Marginalization Procedure.

We now address the question, how such edges can be identified. First, notice that given a DMAG $M$ with an indirect edge $X \to Y$ there always exists a DAG extension $G$ which contains the edge $X \to Y$ (simply include the edge $X \to Y$ in $G$). Thus, it is impossible to identify whether a directed edge is indirect by observing only the DMAG. It is only possible to tell whether an edge is possibly-indirect (i.e. if there is a DAG extension in which the edge does not exist) or definitely direct (i.e. if it exists in all DAG extensions). It turns out that for given DMAG $M$, a directed edge $X \to Y$ is possibly-indirect if (a) there is a directed path from $X$

to $Y$ in $M$ and (b) an inducing path is possible between $X$ and $Y$, i.e. there is an extension of $M$ where it could happen. Notice that the edge $X \to Y$ in $M$ is a directed path from $X$ to $Y$ and an inducing path between $X$ and $Y$. Thus we have to check whether there are *non-direct* directed and inducing paths (i.e. without considering the edge $X \to Y$).

It is easy to check whether the first condition holds; since the Marginalization Procedure preserves ancestral relations, one has to check whether there is a non-direct directed path from $X$ to $Y$. On the other hand, it is not obvious how the second condition can be checked when only observing a DMAG. The reason is that not all latent variables are encoded as bi-directed edges in the DMAG. As we saw in the last section however, we can distinguish between pure-causal and possibly-confounded edges, which comes handy here.

We split the problem solution to two separate cases. Let's first consider the case where a directed edge $X \to Y$ in a DMAG $M$ is labeled as possibly-confounded. In this case, a non-direct inducing path trivially exists between $X$ and $Y$ in a DAG extension $G$ of $M$; just include a latent confounder $H$ of $X$ and $Y$ in $G$, which creates the inducing path $X \leftarrow H \to Y$ in $G$ relative to $L = \{H\}$. Thus, the edge is possibly-indirect if there also exists a non-direct directed path from $X$ to $Y$ in $M$, and is definitely direct otherwise. Now we turn our attention to the other case, i.e. to directed edges that are labeled as pure-causal.

**Theorem 3.** *Let $X \to Y$ be a pure-causal edge in a DMAG $M$. Then $X \to Y$ is possibly-indirect if and only if there is a vertex $W$ in $M$ s.t. $X \to W \to Y$ is in $M$ and $W \to Y$ is possibly-confounded.*

*Proof.* **If**. If $X \to W \to Y$ is in $M$ and $W \to Y$ is possibly-confounded, then there is a DAG extension of $M$ where $W$ and $Y$ share a latent confounder $L_{W,Y}$, and thus the path $\langle X, W, L_{W,Y}, Y \rangle$ is inducing relative to $\mathbf{L} = \{L_{W,Y}\}$. Also the path $\langle X, W, Y \rangle$ is directed from $X$ to $Y$. Thus *if $X \to W \to Y$ is in $M$ then $X \to Y$ is possibly-indirect.*

**And only if**. For the opposite direction, assume that $X \to Y$ is possibly-indirect. Then by definition there exists a DAG extension $G$ of $M$ containing a non-direct directed path from $X$ to $Y$ and a non-direct inducing path $p_{ind}$ from $X$ to $Y$ relative to some set of latent variables **L**.

By definition of inducing paths, every non-latent vertex on $p_{ind}$ is a collider on that path and an ancestor of $X$ or $Y$. Since $X$ is an ancestor of $Y$ ($X \to Y$ is in $M$), *every collider on $p_{ind}$ is an ancestor of $Y$* (if it is an ancestor of $X$ it also is an ancestor of $Y$). Thus, it is easy to see that there is an inducing path from every non-latent vertex on $p_{ind}$ to $Y$. As a result $X \to W \to Y$ or $X \leftarrow *W \to Y$ is in $M$ (where $W$ is the vertex which is adjacent to $X$ and on $p_{ind}$). However, $X \leftarrow *W \to Y$ cannot be in $M$ because it implies that $X$ and $W$ share a latent confounder in $G$ (otherwise $p_{ind}$ would not be an inducing path), which contradicts the fact that there is no inducing path from $X$ to $Y$ relative to the set of latent variables that is into $X$, because $X \to Y$ is pure-causal (given) (see Lemma 9 in (Zhang, 2008))

Thus, *if $X \to Y$ is possibly-indirect, then $X \to W \to Y$ is in $M$*. It remains to show that $W \to Y$ is always possibly-confounded. We will show this by contradiction.

Assume that $W \to Y$ is pure-causal. Then there exists a vertex $U$ in $M$, s.t. $U$ is not connected to $Y$ and (a) $U \to W$ is in $M$, or (b) there is a collider path $p_{col}$ in $M$ from $U$ to $W$ where each vertex is a parent of $Y$.

**Case (a)**: By replacing $X$ with $U$ in $p_{ind}$ we create an inducing path from $U$ to $Y$ (each non-latent vertex is a collider and an ancestor of $Y$), thus $U$ would be adjacent to $Y$ in $M$ and $U \to W$ would not be pure-causal, contradicting our assumption.

**Case (b)**: Let $p'_{col}$ be the path resulting by adding a vertex $L_{V_i, V_{i+1}}$ between each subsequent pair of vertices $V_i$ and $V_{i+1}$ in $p_{col}$, where $L_{V_i, V_{i+1}}$ is a latent parent of $V_i$ and $V_{i+1}$ which replaces the bi-directed edge $V_i \leftrightarrow V_{i+1}$ in $G$. Then there is an inducing path $p'_{ind} = \langle p'_{col}, p_{ind} \setminus X \rangle$, since (a) every non-latent vertex on $p'_{col}$ is a collider, (b) every non-latent vertex on $p_{ind}$ is a collider, (c) $W$ is a collider (the last

vertex on $p'_{col}$ and the first vertex on $p_{ind} \setminus X$ are into $W$), (d) every vertex on $p'_{col}$ is a parent of $Y$, and (d) every vertex on $p_{ind} \setminus X$ is an ancestor of $Y$, thus $U$ and $Y$ would be adjacent in $M$ and $U \to W$ would not be pure-causal, contradicting our assumption.

Thus, *if $X \to Y$ is possibly-indirect, then $X \to W \to Y$ is in $M$ and $W \to Y$ is possibly-confounded*. □

**Theorem 4.** *The time complexity of labeling all directed edges in a MAG $\mathcal{M}$ as direct-causal or possibly-indirect, is $O(n \cdot m + n^2)$.*

*Proof.* As a first step, we have to label each directed edge as pure-causal or possibly-confounded. This takes $O(n \cdot m)$ time (Theorem 2). In order to avoid unnecessary computations we pre-compute the ancestral relations (transitive closure) of $\mathcal{M}$, by ignoring bi-directed edges. This can be trivially done in $O(n \cdot m + n^2)$ time. To label all possibly-confounded edges, we have to consider each edge once, and since all ancestral relations are pre-computed the total time is $O(m)$. Labeling all pure-causal edges takes $O(n \cdot m)$ time, since for each edge $X \to Y$ we have to consider each node at most once (to check if it is connected to $X$ and $Y$), and we have to check whether the remaining conditions are satisfied (see Theorem 3), which can be done in constant time. The total time complexity is $O(n \cdot m + n^2)$. □

## 6 Quantifying the Frequency of the Pure-Causal and Direct-Causal Edges

**Experimental Setup**. We fixed the number of vertices to $n = 40$ and repeated the following experiment 500 times: (a) We varied the edge density $d$ from 0.025 to 1, with a step size of 0.025. For $k$ nodes, the number of edges $m$ is $m = round(d \cdot k \cdot \frac{k-1}{2})$. (b) We varied the number of latent variables $l$ (not necessarily latent confounders) from 0 to 20. (c) We generated uniformly at random a DAG $G$ with $n + l$ nodes and $m$ edges (defined above). The set of latent variables **L** was generated by randomly picking $l$ vertices. A
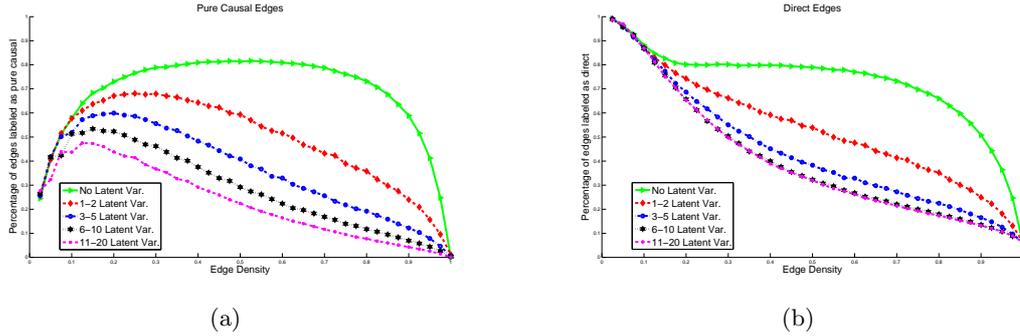
Figure 3: (a) Frequencies of pure-causal and direct-causal edges in randomly generated DMAGs. A relatively high percentage of directed edges can be labeled as pure-causal for graphs with intermediate density. (b) Similarly, a high percentage of directed edges can be labeled as direct-causal for low and medium graph densities.

DMAG $M$ was generated given $G$ and $\mathbf{L}$ using the Marginalization Procedure. Thus, in each DMAG we observe 40 variables, whereas the number of latent variables varies between 0 and 20. (d) We identified all pure causal and direct edges for $M$. The total number of runs is $500 \cdot 40(\text{density}) \cdot 21(\text{latent variables}) = 420000$.

**Results**. The results are shown in Figure 3(a) and Figure 3(b). The x-axes are the edge density (both directed and bi-directed edges) of the resulting DMAG after marginalization of the generating DAG. The y-axes are the percentage of *directed edges* in the DMAG that are labeled as pure-causal or direct-causal respectively. The lines in the plots have been produced by a moving average of a window with range 0.01 and a step of 0.025 to smooth out the results.

There are several observations to make. A relatively high percentage of directed edges can be labeled as pure-causal for graphs with intermediate density (Figure 3(a)). The plots follow a bell shape with the percentage of pure-causal edges dropping to 0 for complete graphs. The percentage of pure-causal edges decreases as the number of latent variables increases. Similarly, a high percentage of directed edges can be labeled as direct-causal for low and medium graph densities (Figure 3(b)). The percentage of direct-causal edges decreases monotonically with increased graph density.

## 7 Conclusion

The causal interpretation of the directed edges $X \rightarrow Y$ in a Maximal Ancestral Graph is not straight-forward and could be misleading to an unfamiliar practitioner. We present efficient algorithms for the theory present by (Zhang, 2008) and extensions that can label directed edges to those that cannot be confounded (pure-causal) or those that cannot be mediated by other observed variables (direct-causal). The algorithms could facilitate interpretation and understanding of a causal graph when latent confounding variables are admitted. The percentage of pure-causal and direct-causal edges depends on the density of the causal graph.

## References

Z. Cai, M. Kuroki, J. Pearl, and J. Tian. 2008. Bounds on Direct Effects in the Presence of Confounded Intermediate Variables. *Biometrics*, 64:695–701.

J. Pearl. 2000. *Causality, Models, Reasoning, and Inference*. Cambridge University Press, Cambridge, U.K.

J. Peters, J. M. Mooij, D. Janzing, and B. Schölkopf. 2012. Identifiability of causal graphs using functional models. *CoRR*, abs/1202.3757.

Th. Richardson and P. Spirtes. 2002. Ancestral graph Markov models. *Annals of Statistics*, 30(4):962–1030.

P. Spirtes, C. Glymour, and R. Scheines. 2000. *Causation, Prediction, and Search*. MIT Press, Cambridge, MA, 2nd edition.

J. Zhang. 2008. Causal reasoning with ancestral graphs. *J. Mach. Learn. Res.*, 9:1437–1474.