

PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a publisher's version.

For additional information about this publication click this link.

<http://hdl.handle.net/2066/204510>

Please be advised that this information was generated on 2019-12-10 and may be subject to change.



Causal Inference by String Diagram Surgery

Bart Jacobs¹, Aleks Kissinger¹, and Fabio Zanasi²(✉)

¹ Radboud University, Nijmegen, The Netherlands

² University College London, London, UK

f.zanasi@ucl.ac.uk

Abstract. Extracting causal relationships from observed correlations is a growing area in probabilistic reasoning, originating with the seminal work of Pearl and others from the early 1990s. This paper develops a new, categorically oriented view based on a clear distinction between syntax (string diagrams) and semantics (stochastic matrices), connected via interpretations as structure-preserving functors.

A key notion in the identification of causal effects is that of an intervention, whereby a variable is forcefully set to a particular value independent of any prior dependencies. We represent the effect of such an intervention as an endofunctor which performs ‘string diagram surgery’ within the syntactic category of string diagrams. This diagram surgery in turn yields a new, interventional distribution via the interpretation functor. While in general there is no way to compute interventional distributions purely from observed data, we show that this is possible in certain special cases using a calculational tool called comb disintegration.

We showcase this technique on a well-known example, predicting the causal effect of smoking on cancer in the presence of a confounding common cause. We then conclude by showing that this technique provides simple sufficient conditions for computing interventions which apply to a wide variety of situations considered in the causal inference literature.

Keywords: Causality · String diagrams · Probabilistic reasoning

1 Introduction

An important conceptual tool for distinguishing correlation from causation is the possibility of *intervention*. For example, a randomised drug trial attempts to destroy any confounding ‘common cause’ explanation for correlations between drug use and recovery by randomly assigning a patient to the control or treatment group, independent of any background factors. In an ideal setting, the observed correlations of such a trial will reflect genuine causal influence. Unfortunately, it is not always possible (or ethical) to ascertain causal effects by means of actual interventions. For instance, one is unlikely to get approval to run a clinical trial on whether smoking causes cancer by randomly assigning 50% of the

patients to smoke, and waiting a bit to see who gets cancer. However, in certain situations it is possible to predict the effect of such a hypothetical intervention from purely observational data.

In this paper, we will focus on the problem of *causal identifiability*. For this problem, we are given observational data as a joint distribution on a set of variables and we are furthermore provided with a *causal structure* associated with those variables. This structure, which typically takes the form of a directed acyclic graph or some variation thereof, tells us which variables can in principle have a causal influence on others. The problem then becomes whether we can measure how strong those causal influences are, by means of computing an *interventional* distribution. That is, can we ascertain what would have happened if a (hypothetical) intervention had occurred?

Over the past 3 decades, a great deal of work has been done in identifying necessary and sufficient conditions for causal identifiability in various special cases, starting with very specific notions such as the *back-door* and *front-door* criteria [20] and progressing to more general necessary and sufficient conditions for causal identifiability based on the **do**-calculus [11], or combinatoric concepts such as confounded components in semi-Markovian models [25, 26].

This style of causal reasoning relies crucially on a delicate interplay between syntax and semantics, which is often not made explicit in the literature. The syntactic object of interest is the causal structure (e.g. a causal graph), which captures something about our understanding of the world, and the mechanisms which gave rise to some observed phenomena. The semantic object of interest is the data: joint and conditional probability distributions on some variables. Fixing a causal structure entails certain constraints on which probability distributions can arise, hence it is natural to see distributions satisfying those constraints as models of the syntax.

In this paper, we make this interplay precise using functorial semantics in the spirit of Lawvere [17], and develop basic syntactic and semantic tools for causal reasoning in this setting. We take as our starting point a functorial presentation of Bayesian networks similar to the one appearing in [7]. The syntactic role is played by string diagrams, which give an intuitive way to represent morphisms of a monoidal category as boxes plugged together by wires. Given a directed acyclic graph (dag) G , we can form a free category Syn_G whose arrows are (formal) string diagrams which represent the causal structure syntactically. Structure-preserving functors from Syn_G to Stoch , the category of stochastic matrices, then correspond exactly to Bayesian networks based on the dag G .

Within this framework, we develop the notion of intervention as an operation of ‘string diagram surgery’. Intuitively, this cuts a string diagram at a certain variable, severing its link to the past. Formally, this is represented as an endofunctor on the syntactic category $\text{cut}_x : \text{Syn}_G \rightarrow \text{Syn}_G$, which propagates through a model $\mathcal{F} : \text{Syn}_G \rightarrow \text{Stoch}$ to send observational probabilities $\mathcal{F}(\omega)$ to interventional probabilities $\mathcal{F}(\text{cut}_x(\omega))$.

The cut_x endofunctor gives us a diagrammatic means of computing interventional distributions given complete knowledge of \mathcal{F} . However, more interestingly,

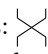
we can sometimes compute interventionals given only partial knowledge of \mathcal{F} , namely some observational data. We show that this can also be done via a technique we call *comb disintegration*, which is a string diagrammatic version of a technique called *c-factorisation* introduced by Tian and Pearl [26]. Our approach generalises disintegration, a calculational tool whereby a joint state on two variables is factored into a single-variable state and a channel, representing the marginal and conditional parts of the distribution, respectively. Disintegration has recently been formulated categorically in [5] and using string diagrams in [4]. We take the latter as a starting point, but instead consider a factorisation of a three-variable state into a channel and a *comb*. The latter is a special kind of map which allows inputs and outputs to be interleaved. They were originally studied in the context of quantum communication protocols, seen as games [8], but have recently been used extensively in the study of causally-ordered quantum [3, 21] and generalised [15] processes. While originally imagined for quantum processes, the categorical formulation given in [15] makes sense in both the classical case (Stoch) and the quantum. Much like Tian and Pearl’s technique, comb factorisation allows one to characterise when the confounding parts of a causal structure are suitably isolated from each other, then exploit that isolation to perform the concrete calculation of interventional distributions.

However, unlike in the traditional formulation, the syntactic and semantic aspects of causal identifiability within our framework exactly mirror one-another. Namely, we can give conditions for causal identifiability in terms of factorisation a morphism in Syn_G , whereas the actual concrete computation of the interventional distribution involves factorisation of its interpretation in Stoch. Thanks to the functorial semantics, the former immediately implies the latter.

To introduce the framework, we make use of a running example taken from Pearl’s book [20]: identifying the causal effect of smoking on cancer with the help of an auxiliary variable (the presence of tar in the lungs). After providing some preliminaries on stochastic matrices and the functorial presentation of Bayesian networks in Sects. 2 and 3, we introduce the smoking example in Sect. 4. In Sect. 5 we formalise the notion of intervention as string diagram surgery, and in Sect. 6 we introduce the combs and prove our main calculational result: the existence and uniqueness of comb factorisations. In Sect. 7, we show how to apply this theorem in computing the interventional distribution in the smoking example, and in 8, we show how this theorem can be applied in a more general case which captures (and slightly generalises) the conditions given in [26]. In Sect. 9, we conclude and describe several avenues of future work.

2 Stochastic Matrices and Conditional Probabilities

Symmetric monoidal categories (SMCs) give a very general setting for studying processes which can be composed in sequence (via the usual categorical composition \circ) and in parallel (via the monoidal composition \otimes). Throughout this paper, we will use *string diagram* notation [24] for depicting composition of morphisms in an SMC. In this notation, morphisms are depicted as boxes with labelled input

and output wires, composition \circ as ‘plugging’ boxes together, and the monoidal product \otimes as placing boxes side-by-side. Identity morphisms are depicted simply as a wire and the unit I of \otimes as the empty diagram. The ‘symmetric’ part of the structure consists of symmetry morphisms, which enable us to permute inputs and outputs arbitrarily. We depict these as wire-crossings: . Morphisms whose domain is I are called *states*, and they will play a special role throughout this paper.

A monoidal category of prime interest in this paper is **Stoch**, whose objects are finite sets and morphisms $\mathbf{f} : A \rightarrow B$ are $|B| \times |A|$ dimensional stochastic matrices. That is, they are matrices of positive numbers (including 0) whose columns each sum to 1:

$$\mathbf{f} = \{\mathbf{f}_i^j \in \mathbb{R}^+ \mid i \in A, j \in B\} \quad \text{with} \quad \sum_j \mathbf{f}_i^j = 1, \text{ for all } i.$$

Note we adopt the physicists convention of writing row indices as superscripts and column indices as subscripts. Stochastic matrices are of interest for probabilistic reasoning, because they exactly capture the data of a conditional probability distribution. That is, if we take $A := \{1, \dots, m\}$ and $B := \{1, \dots, n\}$, conditional probabilities naturally arrange themselves into a stochastic matrix:

$$\mathbf{f}_i^j := P(B = j \mid A = i) \quad \rightsquigarrow \quad \mathbf{f} = \begin{pmatrix} P(B = 1 \mid A = 1) & \cdots & P(B = 1 \mid A = m) \\ \vdots & \ddots & \vdots \\ P(B = n \mid A = 1) & \cdots & P(B = n \mid A = m) \end{pmatrix}$$

States, i.e. stochastic matrices from a trivial input $I := \{*\}$, are (non-conditional) probability distributions, represented as column vectors. There is only one stochastic matrix with trivial output: the row vector consisting only of 1’s. The latter, with notation \blackuparrow as on the right, will play a special role in this paper (see (1) below).

Composition of stochastic matrices is matrix multiplication. In terms of conditional probabilities, that is multiplication followed by marginalization over the shared variable: $\sum_B P(C \mid B)P(B \mid A)$. Identities are thus given by identity matrices, which we will often express in terms of the Kronecker delta function δ_i^j .

The monoidal product \otimes in **Stoch** is the cartesian product on objects, and Kronecker product of matrices: $(\mathbf{f} \otimes \mathbf{g})_{(i,j)}^{(k,l)} := \mathbf{f}_i^k \mathbf{g}_j^l$. We will typically omit parentheses and commas in the indices, writing e.g. \mathbf{h}_{ij}^{kl} instead of $\mathbf{h}_{(i,j)}^{(k,l)}$ for an arbitrary matrix entry of $\mathbf{h} : A \otimes B \rightarrow C \otimes D$. In terms of conditional probabilities, the Kronecker product corresponds to taking product distributions. That is, if \mathbf{f} represents the conditional probabilities $P(B \mid A)$ and \mathbf{g} the probabilities $P(D \mid C)$, then $\mathbf{f} \otimes \mathbf{g}$ represents $P(B \mid A)P(D \mid C)$. **Stoch** also comes with a natural choice of ‘swap’ matrices $\sigma : A \otimes B \rightarrow B \otimes A$ given by $\sigma_{ij}^{kl} := \delta_i^l \delta_j^k$, making it into a symmetric monoidal category. Every object A in **Stoch** has three other pieces of structure which will play a key role in our formulation of Bayesian networks and interventions: the *copy* map, the *discarding* map, and the *uniform state*:

$$\left(\begin{array}{c} \diagdown \\ \bullet \\ \diagup \end{array} \right)_i^{jk} := \delta_i^j \delta_i^k \qquad \left(\bullet \right)_i := 1 \qquad \left(\begin{array}{c} \downarrow \end{array} \right)^i := \frac{1}{|A|} \qquad (1)$$

Abstractly, this provides **Stoch** with the structure of a *CDU category*.

Definition 2.1. A CDU category (for *copy, discard, uniform*) is a symmetric monoidal category (\mathbb{C}, \otimes, I) where each object A has a copy map $\blacktriangledown : A \rightarrow A \otimes A$, a discarding map $\blacklozenge : A \rightarrow I$, and a uniform state $\blacktriangledown : I \rightarrow A$ satisfying the following equations:

CDU functors are symmetric monoidal functors between CDU categories preserving copy maps, discard maps and uniform states.

We assume that the CDU structure on I is trivial and the CDU structure on $A \otimes B$ is constructed in the obvious way from the structure on A and B . We also use the first equation in (2) to justify writing ‘copy’ maps with arbitrarily many output wires: \blacktriangledown .

Similar to [2], we can form the free CDU category $\text{FreeCDU}(X, \Sigma)$ over a pair (X, Σ) of a generating set of objects X and a generating set Σ of typed morphisms $f : u \rightarrow w$, with $u, w \in X^*$ as follows. The category $\text{FreeCDU}(X, \Sigma)$ has X^* as set of objects, and morphisms the string diagrams constructed from the elements of Σ and maps $\blacktriangledown : x \rightarrow x \otimes x$, $\blacklozenge : x \rightarrow I$ and $\blacktriangledown : I \rightarrow x$ for each $x \in X$, taken modulo the equations (2).

Lemma 2.2. *Stoch* is a CDU category, with CDU structure defined as in (1).

An important feature of **Stoch** is that $I = \{\star\}$ is the final object, with $\blacklozenge : B \rightarrow I$ the map provided by the universal property, for any set B . This yields Eq. (3) on the right, for any $f : A \rightarrow B$, justifying the name “discarding map” for \blacklozenge .

We conclude by recording another significant feature of **Stoch**: *disintegration* [4, 5]. In probability theory, this is the mechanism of factoring a joint probability distribution $P(AB)$ as a product of the first marginal $P(A)$ and a conditional distribution $P(B|A)$. We recall from [4] the string diagrammatic rendition of this process. We say that a morphism $f : X \rightarrow Y$ in **Stoch** has *full support* if, as a stochastic matrix, it has no zero entries. When f is a state, it is a standard result that full support ensures uniqueness of disintegrations of f .

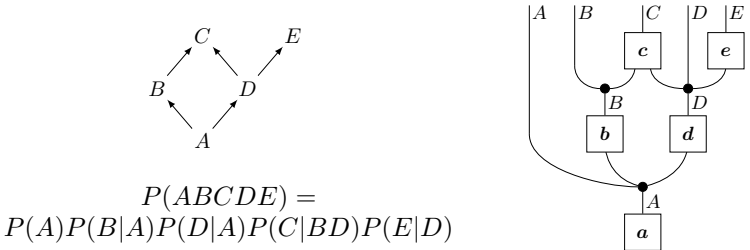
Proposition 2.3 (Disintegration). *For any state $\omega : I \rightarrow A \otimes B$ with full support, there exists unique morphisms $a : I \rightarrow A, b : A \rightarrow B$ such that:*

Note that Eq. (3) and the CDU rules immediately imply that the unique $\mathbf{a}: I \rightarrow A$ in Proposition 2.3 is the marginal of ω onto A : $\begin{array}{|c|} \hline A \\ \hline \omega \\ \hline \end{array}$.

3 Bayesian Networks as String Diagrams

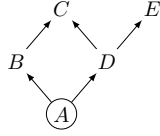
Bayesian networks are a widely-used tool in probabilistic reasoning. They give a succinct representation of conditional (in)dependencies between variables as a directed acyclic graph. Traditionally, a Bayesian network on a set of variables A, B, C, \dots is defined as a directed acyclic graph (dag) G , an assignment of sets to each of the nodes $V_G := \{A, B, C, \dots\}$ of G and a joint probability distribution over those variables which factorises as $P(V_G) = \prod_{A \in V_G} P(A | \text{Pa}(A))$ where $\text{Pa}(A)$ is the set of parents of A in G . Any joint distribution that factorises this way is said to satisfy the *global Markov property* with respect to the dag G . Alternatively, a Bayesian network can be seen as a dag equipped with a set of conditional probabilities $\{P(A | \text{Pa}(A)) \mid A \in V_G\}$ which can be combined to form the joint state. Thanks to disintegration, these two perspectives are equivalent.

Much like in the case of disintegration in the previous section, Bayesian networks have a neat categorical description as string diagrams in the category *Stoch* [7, 13, 14]. For example, here is a Bayesian network in its traditional depiction as a dag with an associated joint distribution over its vertices, and as a string diagram in *Stoch*:

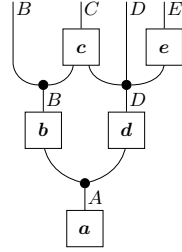


In the string diagram above, the stochastic matrix $\mathbf{a}: I \rightarrow A$ contains the probabilities $P(A)$, $\mathbf{b}: B \rightarrow A$ contains the conditional probabilities $P(B|A)$, $\mathbf{c}: B \otimes D \rightarrow C$ contains $P(C|BD)$, and so on. The entire diagram is then equal to a state $\omega: I \rightarrow A \otimes B \otimes C \otimes D \otimes E$ which represents $P(ABCDE)$.

Note the dag and the diagram above look similar in structure. The main difference is the use of copy maps to make each variable (even those that are not leaves of the dag, A, B and D) an output of the overall diagram. This corresponds to a variable being *observed*. We can also consider Bayesian networks with *latent* variables, which do not appear in the joint distribution due to marginalisation. Continuing the example above, making A into a latent variable yields the following depiction as a string diagram:



$$P(BCDE) = \sum_A P(A)P(B|A)P(D|A)P(C|BD)P(E|D)$$



In general, a Bayesian network (with possible latent variables), is a string diagram in **Stoch** that (1) only has outputs and (2) consists only of copy maps and boxes which each have exactly one output.

By ‘a string diagram in **Stoch**’, we mean not only the stochastic matrix itself, but also its decomposition into components. We can formalise exactly what we mean by taking a perspective on Bayesian networks which draws inspiration from Lawvere’s functorial semantics of algebraic theories [16]. In this perspective, which elaborates on [7, Ch. 4], we maintain a conceptual distinction between the purely syntactic object (the diagram) and its probabilistic interpretation.

Starting from a dag $G = (V_G, E_G)$, we construct a free CDU category Syn_G which provides the syntax of causal structures labelled by G . The objects of Syn_G are generated by the vertices of G , whereas the morphisms are generated by the following signature:

$$\Sigma_G = \left\{ \left[\begin{array}{c} | \\ \text{A} \\ | \\ \boxed{a} \\ | \\ \text{B}_1 \cdots \text{B}_k \end{array} \right] \mid A \in V_G \text{ with parents } B_1, \dots, B_k \in V_G \right\}$$

Then $\text{Syn}_G := \text{FreeCDU}(V_G, \Sigma_G)$.¹ The following result establishes that models (à la Lawvere) of Syn_G coincide with G -based Bayesian networks.

Proposition 3.1. *There is a 1-1 correspondence between Bayesian networks based on the dag G and CDU functors of type $\text{Syn}_G \rightarrow \text{Stoch}$.*

We refer to [12] for a proof. This proposition justifies the following definition of a category BN_G of G -based Bayesian networks: objects are CDU functors $\text{Syn}_G \rightarrow \text{Stoch}$ and arrows are monoidal natural transformations between them.

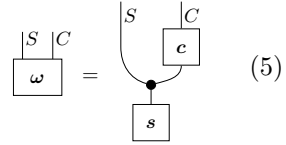
4 Towards Causal Inference: The Smoking Scenario

We will motivate our approach to causal inference via a classic example, inspired by the one given in the Pearl’s book [20]. Imagine a dispute between a scientist and a tobacco company. The scientist claims that smoking causes cancer. As a source of evidence, the scientist cites a joint probability distribution ω over variables S for smoking and C for cancer, which disintegrates as in (5) below,

¹ Note that E_G is implicitly used in the construction of Syn_G : the edges of G determine the parents of a vertex, and hence the input types of the symbols in Σ_G .

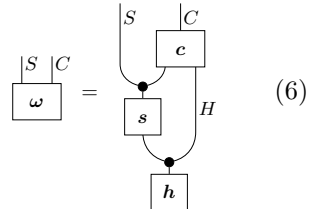
with matrix $\mathbf{c} = \begin{pmatrix} 0.9 & 0.7 \\ 0.1 & 0.3 \end{pmatrix}$. Inspecting this $\mathbf{c} : S \rightarrow C$, the scientist notes that the probability of getting cancer for smokers (0.3) is three times as high as for non-smokers (0.1). Hence, the scientist claims that smoking has a significant causal effect on cancer.

An important thing to stress here is that the scientist draws this conclusion using not only the observational data ω but also from an assumed *causal structure* which gave rise to that data, as captured in the diagram in Eq. (5). That is, rather than treating diagram (5) simply as a calculation on the observational data, it can also be

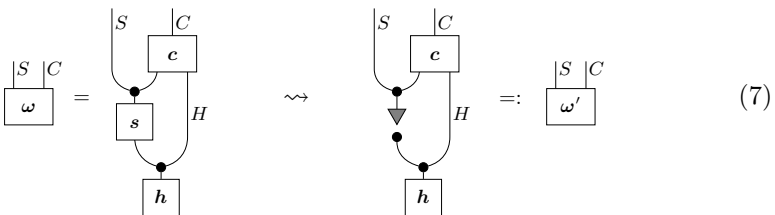


treated as an assumption about the actual, physical mechanism that gave rise to that data. Namely, this diagram encompasses the assumption that there is some prior propensity for people to smoke captured by $s : I \rightarrow S$, which is both observed and fed into some other process $c : S \rightarrow C$ whereby an individuals choice to smoke determines whether or not they get cancer.

The tobacco company, in turn, says that the scientists’ assumptions about the provenance of this data are too strong. While they concede that *in principle* it is possible for smoking to have some influence on cancer, the scientist should allow for the possibility that there is some latent common cause (e.g. genetic conditions, stressful work environment, etc.) which leads people both to smoke and get cancer. Hence, says the tobacco company, a ‘more honest’ causal structure to ascribe to the data ω is (6). This structure then allows for either party to be correct. If the scientist is right, the output of $c : S \otimes H \rightarrow C$ depends mostly on its first input, i.e. the causal path from smoking to cancer. If the tabacco company is right, then c depends very little on its first input, and the correlation between S and C can be explained almost entirely from the hidden common cause.



So, who is right after all? Just from the observed distribution ω , it is impossible to tell. So, the scientist proposes a clinical trial, in which patients are randomly required to smoke or not to smoke. We can model this situation by replacing s in (6) with a process that ignores its inputs and outputs the uniform state. Graphically, this looks like ‘cutting’ the link s between H and S :

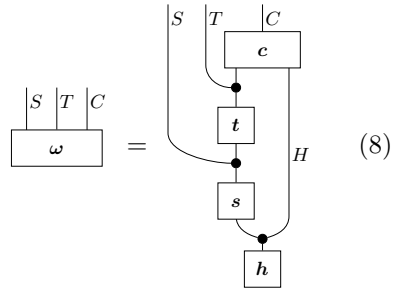


This captures the fact that variable S is now randomised and no longer dependent on any background factors. This new distribution ω' represents the data

the scientist would have obtained had they run the trial. That is, it gives the results of an *intervention* at s . If this ω' *still* shows a strong correlation between smoking and cancer, one can conclude that smoking indeed causes cancer even when we assume the weaker causal structure (6).

Unsurprisingly, the scientist fails to get ethical approval to run the trial, and hence has only the observational data ω to work with. Given that the scientist only knows ω (and not c and h), there is no way to compute ω' in this case. However, a key insight of statistical causal inference is that sometimes it *is* possible to compute interventional distributions from observational ones. Continuing the smoking example, suppose the scientist proposes the following revision to the causal structure: they posit a structure (8) that includes a third observed variable (the presence of T of tar in the lungs), which completely mediates the causal effect of smoking on cancer.

As with our simpler structure, the diagram (8) contains some assumptions about the provenance of the data ω . In particular, by omitting wires, we are asserting there is no *direct* causal link between certain variables. The absence of an H -labelled input to t says there is no direct causal link from H to T (only mediated by S), and the absence of an S -labelled input wire into c captures that there is no direct causal link from S to C (only mediated by T). In the traditional approach to causal inference, such relationships are typically captured by a graph-theoretic property called *d-separation* on the dag associated with the causal structure.



We can again imagine intervening at S by replacing $s : H \rightarrow S$ by $\downarrow \circ \uparrow$. Again, this ‘cutting’ of the diagram will result in a new interventional distribution ω' . However, unlike before, it *is* possible to compute this distribution from the observational distribution ω .

However, in order to do that, we first need to develop the appropriate categorical framework. In Sect. 5, we will model ‘cutting’ as a functor. In 6, we will introduce a generalisation of disintegration, which we call *comb disintegration*. These tools will enable us to compute ω' for ω , in Sect. 7.

5 Interventional Distributions as Diagram Surgery

The goal of this section is to define the ‘cut’ operation in (7) as an endofunctor on the category of Bayesian networks. First, we observe that such an operation exclusively concerns the string diagram part of a Bayesian network: following the functorial semantics given in Sect. 3, it is thus appropriate to define cut as an endofunctor on Syn_G , for a given dag G .

Definition 5.1. For a fixed node $A \in V_G$ in a graph G , let $\text{cut}_A: \text{Syn}_G \rightarrow \text{Syn}_G$ be the CDU functor freely obtained by the following action on the generators (V_G, Σ_G) of Syn_G :

- For each object $B \in V_G$, $\text{cut}_A(B) = B$.
- $\text{cut}_A\left(\begin{array}{c} \boxed{a} \\ |_{B_1} \cdots |_{B_k} \end{array}\right) = \begin{array}{c} \downarrow^A \\ \bullet \\ |_{B_1} \cdots |_{B_k} \end{array}$ and $\text{cut}_A\left(\begin{array}{c} \boxed{b} \\ |_{C_1} \cdots |_{C_j} \end{array}\right) = \begin{array}{c} \boxed{b} \\ |_{C_1} \cdots |_{C_j} \end{array}$ for any other $\begin{array}{c} \boxed{b} \\ |_{C_1} \cdots |_{C_j} \end{array} \in \Sigma_G$.

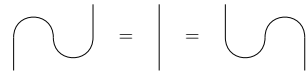
Intuitively, cut_A applied to a string diagram f of Syn_G removes from f each occurrence of a box with output wire of type A .

Proposition 3.1 allows us to “transport” the cutting operation over to Bayesian networks. Given any Bayesian network based on G , let $\mathcal{F}: \text{Syn}_G \rightarrow \text{Stoch}$ be the corresponding CDU functor given by Proposition 3.1. Then, we can define its A -cutting as the Bayesian network identified by the CDU functor $\mathcal{F} \circ \text{cut}_A$. This yields an (idempotent) endofunctor $\text{Cut}_A: \text{BN}_G \rightarrow \text{BN}_G$.

6 The Comb Factorisation

Thanks to the developments of Sect. 5, we can understand the transition from left to right in (7) as the application of the functor Cut_s applied to the ‘Smoking’ node S . The next step is being able to actually compute the individual Stoch -morphisms appearing in (8), to give an answer to the causality question.

In order to do that, we want to work in a setting where $t: S \rightarrow T$ can be isolated and ‘extracted’ from (8). What is left behind is a stochastic matrix with a ‘hole’ where t has been



extracted. To define ‘morphisms with holes’, it is convenient to pass from SMCs to compact closed categories (see e.g. [24]). Stoch is not itself compact closed, but it embeds into $\text{Mat}(\mathbb{R}^+)$, whose morphisms are *all* matrices over positive numbers. $\text{Mat}(\mathbb{R}^+)$ has a (self-dual) compact closed structure; that means, for any set A there is a ‘cap’ $\cap: A \otimes A \rightarrow I$ and a ‘cup’ $\cup: I \rightarrow A \otimes A$, which satisfy the ‘yanking’ equations on the right. As matrices, caps and cups are defined by $\cap_{ij} = \cup^{ij} = \delta_i^j$. Intuitively, they amount to ‘bent’ identity wires. Another aspect of $\text{Mat}(\mathbb{R}^+)$ that is useful to recall is the following handy characterisation of the subcategory Stoch .

Lemma 6.1. A morphism $f: A \rightarrow B$ in $\text{Mat}(\mathbb{R}^+)$ is a stochastic matrix (thus a morphism of Stoch) if and only if (3) holds.

A suitable notion of ‘stochastic map with a hole’ is provided by a *comb*. These structures originate in the study of certain kinds of quantum channels [3].

Definition 6.2. A 2-comb in Stoch is a morphism $f: A_1 \otimes A_2 \rightarrow B_1 \otimes B_2$ satisfying, for some other morphism $f': A_1 \rightarrow B_1$,

$$\begin{array}{c} \begin{array}{c} B_1 | \\ \boxed{f} \\ A_1 | \quad A_2 \end{array} \begin{array}{c} \bullet \\ |_{B_2} \end{array} = \begin{array}{c} \begin{array}{c} B_1 | \\ \boxed{f'} \\ A_1 | \end{array} \begin{array}{c} \bullet \\ |_{A_2} \end{array} \end{array} \tag{9}$$

This definition extends inductively to n -combs, where we require that discarding the rightmost output yields $\mathbf{f}' \otimes \uparrow$, for some $(n - 1)$ -comb \mathbf{f}' . However, for our purposes, restricting to 2-combs will suffice.

The intuition behind condition (9) is that the contribution from input A_2 is only visible via output B_2 . Thus, if we discard B_2 we may as well discard A_2 . In other words, the input/output pair A_2, B_2 happen ‘after’ the pair A_1, B_1 . Hence, it is typical to depict 2-combs in the shape of a (hair) comb, with 2 ‘teeth’, as in (10) below:

While combs themselves live in Stoch , $\text{Mat}(\mathbb{R}^+)$ accommodates a second-order reading of the transition \rightsquigarrow in (10): we can treat \mathbf{f} as a map which expects as input a map $\mathbf{g} : B_1 \rightarrow A_2$ and produces as output a map of type $A_1 \rightarrow B_2$. Plugging $\mathbf{g} : B_1 \rightarrow A_2$ into the 2-comb can be formally defined in $\text{Mat}(\mathbb{R}^+)$ by composing \mathbf{f} and \mathbf{g} in the usual way, then feeding the output of \mathbf{g} into the second input of \mathbf{f} , using caps and cups, as in (11).

Importantly, for generic \mathbf{f} and \mathbf{g} of Stoch , there is no guarantee that forming the composite (11) in $\text{Mat}(\mathbb{R}^+)$ yields a valid Stoch -morphism, i.e. a morphism satisfying the finality Eq. (3). However, if \mathbf{f} is a 2-comb and \mathbf{g} is a Stoch -morphism, Eq. (9) enables a discarding map plugged into the output B_2 in (11) to ‘fall through’ the right side of \mathbf{f} , which guarantees that the composed map satisfies the finality equation for discarding. See [12, § ??] for the explicit diagram calculation.

With the concept of 2-combs in hand, we can state our factorisation result.

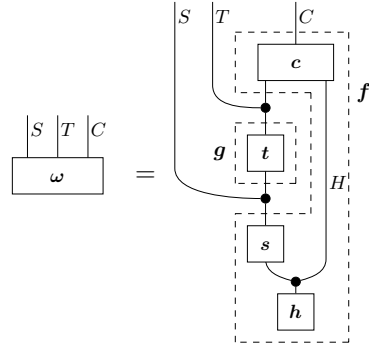
Theorem 6.3. *For any state $\omega : I \rightarrow A \otimes B \otimes C$ of Stoch with full support, there exists a unique 2-comb $\mathbf{f} : B \rightarrow A \otimes C$ and stochastic matrix $\mathbf{g} : A \rightarrow B$ such that, in $\text{Mat}(\mathbb{R}^+)$:*

Proof. The construction of \mathbf{f} and \mathbf{g} mimics the one of c-factors in [26], using string diagrams and (diagrammatic) disintegration. We first use ω to construct maps $\mathbf{a} : I \rightarrow A, \mathbf{b} : A \rightarrow B, \mathbf{c} : A \otimes B \rightarrow C$, then construct \mathbf{f} using \mathbf{a} and \mathbf{c} and construct \mathbf{g} using \mathbf{b} . For the full proof, including uniqueness, see [12]. \square

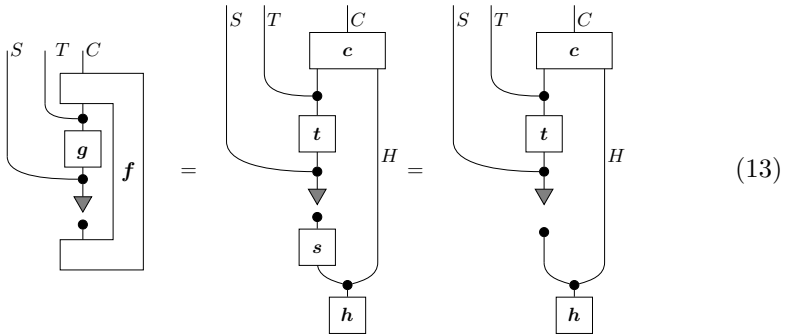
Note that Theorem 6.3 generalises the normal disintegration property given in Proposition 2.3. The latter is recovered by taking $A := I$ (or $C := I$) above.

7 Returning to the Smoking Scenario

We now return to the smoking scenario of Sect. 4. There, we concluded by claiming that the introduction of an intermediate variable T to the observational distribution $\omega : I \rightarrow S \otimes T \otimes C$ would enable us to calculate the interventional distribution. That is, we can calculate $\omega' = \mathcal{F}(\text{cut}_s(\omega))$ from $\omega := \mathcal{F}(\omega)$. Thanks to Theorem 6.3, we are now able to perform that calculation. We first observe that our assumed causal structure for ω fits the form of Theorem 6.3, where g is t and f is a 2-comb containing everything else, as in the diagram on the side.



Hence, f and g are computable from ω . If we plug them back together as in (12), we will get ω back. However, if we insert a ‘cut’ between f and g :



we obtain $\omega' = \mathcal{F}(\text{cut}_s(\omega))$.

We now consider a concrete example. Fix interpretations $S = T = C = \{0, 1\}$ and let $\omega : I \rightarrow S \otimes T \otimes C$ be the stochastic matrix:

$$\omega := \begin{pmatrix} 0.5 \\ 0.1 \\ 0.01 \\ 0.02 \\ 0.1 \\ 0.05 \\ 0.02 \\ 0.2 \end{pmatrix} \begin{matrix} \leftarrow P(S = 0, T = 0, C = 0) \\ \leftarrow P(S = 0, T = 0, C = 1) \\ \leftarrow P(S = 0, T = 1, C = 0) \\ \leftarrow P(S = 0, T = 1, C = 1) \\ \leftarrow P(S = 1, T = 0, C = 0) \\ \leftarrow P(S = 1, T = 0, C = 1) \\ \leftarrow P(S = 1, T = 1, C = 0) \\ \leftarrow P(S = 1, T = 1, C = 1) \end{matrix}$$

Now, disintegrating ω :

$$\omega = c \circ s \quad \text{gives} \quad c \approx \begin{pmatrix} 0.81 & 0.32 \\ 0.19 & 0.68 \end{pmatrix}$$

The bottom-left element of c is $P(C = 1|S = 0)$, whereas the bottom-right is $P(C = 1|S = 1)$, so this suggests that patients are ≈ 3.5 times as likely to get cancer if they smoke (68% vs. 19%). However, comb-disintegrating ω using Theorem 6.3 gives $g: S \rightarrow T$ and a comb $f: T \rightarrow S \otimes C$ with the following stochastic matrices:

$$f \approx \begin{pmatrix} 0.53 & 0.21 \\ 0.11 & 0.42 \\ 0.25 & 0.03 \\ 0.12 & 0.34 \end{pmatrix} \quad g \approx \begin{pmatrix} 0.95 & 0.41 \\ 0.05 & 0.59 \end{pmatrix}$$

Recomposing these with a ‘cut’ in between, as in the left-hand side of (13), gives the interventional distribution $\omega' \approx (0.38, 0.11, 0.01, 0.02, 0.16, 0.05, 0.07, 0.22)$. Disintegrating:

$$\omega' = c' \circ s' \quad \text{gives} \quad c' \approx \begin{pmatrix} 0.75 & 0.46 \\ 0.25 & 0.54 \end{pmatrix}$$

From the interventional distribution, we conclude that, in a (hypothetical) clinical trial, patients are about twice as likely to get cancer if they smoke (54% vs. 25%). So, since $54 < 68$, there was *some* confounding influence between S and C in our observational data, but after removing it via comb disintegration, we see there is still a significant causal link between smoking and cancer.

Note this conclusion depends totally on the particular observational data that we picked. For a different interpretation of ω in *Stoch*, one might conclude that there is *no* causal connection, or even that smoking *decreases* the chance of getting cancer. Interestingly, all three cases can arise even when a naïve analysis of the data shows a strong direct correlation between S and C . To see and/or experiment with these cases, we have provided the Python code² used to perform these calculations. See also [19] for a pedagogical overview of this example (using traditional Bayesian network language) with some sample calculations.

8 The General Case for a Single Intervention

While we applied the comb decomposition to a particular example, this technique applies essentially unmodified to many examples where we intervene at a single variable (called X below) within an arbitrary causal structure.

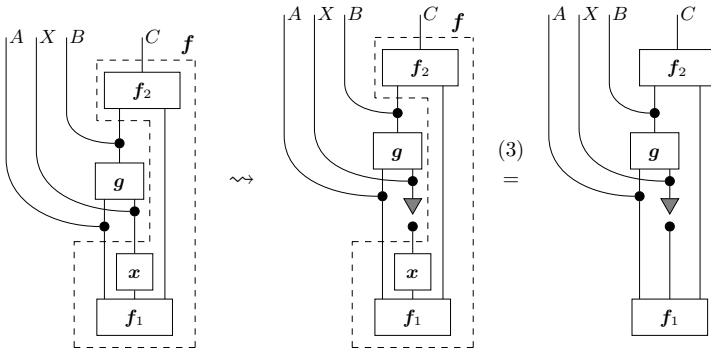
² <https://gist.github.com/akissinger/aec1751792a208253bda491ead587b6>.

Theorem 8.1. *Let G be a dag with a fixed node X that has corresponding generator $x: Y_1 \otimes \dots \otimes Y_n \rightarrow X$ in Syn_G . Then, suppose ω is a morphism in Syn_G of the following form:*

$$\begin{array}{c} \begin{array}{|c|c|c|c|} \hline A & X & B & C \\ \hline \end{array} \\ \hline \omega \end{array} = \begin{array}{c} \begin{array}{|c|} \hline A \\ \hline \end{array} \begin{array}{|c|} \hline X \\ \hline \end{array} \begin{array}{|c|} \hline B \\ \hline \end{array} \begin{array}{|c|} \hline C \\ \hline \end{array} \\ \hline \begin{array}{c} \begin{array}{|c|} \hline f_2 \\ \hline \end{array} \\ \bullet \\ \begin{array}{|c|} \hline g \\ \hline \end{array} \\ \bullet \\ \begin{array}{|c|} \hline x \\ \hline \end{array} \\ \bullet \\ \begin{array}{|c|} \hline f_1 \\ \hline \end{array} \end{array} \end{array} \quad (14)$$

for some morphisms f_1, f_2 and g in Syn_G not containing x as a subdiagram. Then the interventional distribution $\omega' := \mathcal{F}(\text{cut}_x(\omega))$ is computable from the observational distribution $\omega = \mathcal{F}(\omega)$.

Proof. The proof is very close to the example in the previous section. Interpreting ω into Stoch , we get a diagram of stochastic maps, which we can comb-disintegrate, then recompose with $\blacktriangledown \circ \bullet$ to produce the interventional distribution:



The RHS above is then $\mathcal{F}(\text{cut}_x(\omega))$. □

This is general enough to cover several well-known sufficient conditions from the causality literature, including single-variable versions of the so-called *front-door* and *back-door* criteria, as well as the sufficient condition based on confounding paths given by Pearl and Tian [26]. As the latter subsumes the other two, we will say a few words about the relationship between the Pearl/Tian condition and Theorem 8.1. In [26], the authors focus on *semi-Markovian* models, where the only latent variables have exactly two observed children and no parents. Suppose we write $A \leftrightarrow B$ if two observed variables are connected by a latent common cause, then one can characterise *confounding paths* as the transitive closure of \leftrightarrow . They go on to show that the interventional distribution corresponding cutting X is computable whenever there are no confounding paths connecting X to one of its children.

We can compare this to the form of expression ω in Eq. (14). First, note this factorisation implies that all boxes which take X as an input must occur as sub-diagrams of g . Hence, any ‘confounding path’ connecting X to its children would yield at least one (un-copied) wire from f_1 to g , hence it cannot be factored as (14). Conversely, if there are no confounding paths from X to its children, then we can place the boxes involved in any other confounding path either entirely inside of g or entirely outside of g and obtain factorisation (14). Hence, restricting to semi-Markovian models, the no-confounding-path condition from [26] is equivalent to ours. However, Theorem 8.1 is slightly more general: its formulation doesn’t rely on the causal structure ω being semi-Markovian.

9 Conclusion and Future Work

This paper takes a fresh, systematic look at the problem of causal identifiability. By clearly distinguishing syntax (string diagram surgery and identification of comb shapes) and semantics (comb-disintegration of joint states) we obtain a clear methodology for computing interventional distributions, and hence causal effects, from observational data.

A natural next step is moving beyond single-variable interventions to the general case, i.e. situations where we allow interventions on multiple variables which may have some arbitrary causal relationships connecting them. This would mean extending the comb factorisation Theorem 6.3 from a 2-comb and a channel to arbitrary n -combs. This seems to be straightforward, via an inductive extension of the proof of Theorem 6.3. A more substantial direction of future work will be the strengthening of Theorem 8.1 from sufficient conditions for causal identifiability to a full characterisation. Indeed, the related condition based on confounding paths from [26] is a necessary and sufficient condition for computing the interventional distribution on a single variable. Hence, it will be interesting to formalise this necessity proof (and more general versions, e.g. [10]) within our framework and investigate, for example, the extent to which it holds beyond the semi-Markovian case.

While we focus exclusively on the case of taking models in *Stoch* in this paper, the techniques we gave are posed at an abstract level in terms of composition and factorisation. Hence, we are optimistic about their prospects to generalise to other probabilistic (e.g. infinite discrete and continuous variables) and quantum settings. In the latter case, this could provide insights into the emerging field of *quantum causal structures* [6, 9, 18, 22, 23], which attempts in part to replay some of the results coming from statistical causal reasoning, but where quantum processes play a role analogous to stochastic ones. A key difficulty in applying our framework to a category of quantum processes, rather than *Stoch*, is the unavailability of ‘copy’ morphisms due to the quantum no-cloning theorem [27]. However, a recent proposal for the formulation of ‘quantum common causes’ [1] suggests a (partially-defined) analogue to the role played by ‘copy’ in our formulation constructed via multiplication of certain commuting Choi matrices. Hence, it may yet be possible to import results from classical causal reasoning into the quantum case just by changing the category of models.

Acknowledgements. FZ acknowledges support from EPSRC grant EP/R020604/1. AK would like to thank Tom Claassen for useful discussions on causal identification criteria.

References

1. Allen, J.-M.A., Barrett, J., Horsman, D.C., Lee, C.M., Spekkens, R.W.: Quantum common causes and quantum causal models. *Phys. Rev. X* **7**, 031021 (2017)
2. Bonchi, F., Sobociński, P., Zanasi, F.: Deconstructing Lawvere with distributive laws. *J. Log. Algebr. Meth. Program.* **95**, 128–146 (2018)
3. Chiribella, G., D’Ariano, G.M., Perinotti, P.: Quantum circuit architecture. *Phys. Rev. Lett.* **101**, 060401 (2008)
4. Cho, K., Jacobs, B.: Disintegration and Bayesian inversion, both abstractly and concretely (2017). arxiv.org/abs/1709.00322
5. Clerc, F., Danos, V., Dahlqvist, F., Garnier, I.: Pointless learning. In: Esparza, J., Murawski, A.S. (eds.) *FoSSaCS 2017*. LNCS, vol. 10203, pp. 355–369. Springer, Heidelberg (2017). https://doi.org/10.1007/978-3-662-54458-7_21
6. Costa, F., Shrapnel, S.: Quantum causal modelling. *New J. Phys.* **18**(6), 063032 (2016)
7. Fong, B.: Causal theories: a categorical perspective on Bayesian networks. Master’s thesis, University of Oxford (2012). arxiv.org/abs/1301.6201
8. Gutoski, G., Watrous, J.: Toward a general theory of quantum games. In: *Proceedings of the Thirty-Ninth Annual ACM Symposium on Theory of Computing*, pp. 565–574. ACM (2007)
9. Henson, J., Lal, R., Pusey, M.F.: Theory-independent limits on correlations from generalized Bayesian networks. *New J. Phys.* **16**(11), 113043 (2014)
10. Huang, Y., Valtorta, M.: On the completeness of an identifiability algorithm for semi-Markovian models. *Ann. Math. Artif. Intell.* **54**(4), 363–408 (2008)
11. Huang, Y., Valtorta, M.: Pearl’s calculus of intervention is complete. *CoRR*, abs/1206.6831 (2012)
12. Jacobs, B., Kissinger, A., Zanasi, F.: Causal inference by string diagram surgery. *CoRR*, abs/1811.08338 (2018)
13. Jacobs, B., Zanasi, F.: A predicate/state transformer semantics for Bayesian learning. *Electr. Notes Theor. Comput. Sci.* **325**, 185–200 (2016)
14. Jacobs, B., Zanasi, F.: The logical essentials of Bayesian reasoning. *CoRR*, abs/1804.01193 (2018)
15. Kissinger, A., Uijlen, S.: A categorical semantics for causal structure. In: *32nd Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2017, Reykjavik, Iceland, 20–23 June 2017*, pp. 1–12 (2017)
16. Lawvere, F.W.: Ordinal sums and equational doctrines. In: Eckmann, B. (ed.) *Seminar on Triples and Categorical Homology Theory*. LNM, vol. 80, pp. 141–155. Springer, Heidelberg (1969). <https://doi.org/10.1007/BFb0083085>
17. Lawvere, F.W.: Functorial semantics of algebraic theories. *Proc. Natl. Acad. Sci. U.S.A.* **50**(5), 869 (1963)
18. Leifer, M.S., Spekkens, R.W.: Towards a formulation of quantum theory as a causally neutral theory of Bayesian inference. *Phys. Rev. A* **88**, 052130 (2013)
19. Nielsen, M.: If correlation doesn’t imply causation, then what does? <http://www.michaelnielsen.org/ddi/if-correlation-doesnt-imply-causation-then-what-does>. Accessed 15 Nov 2018

20. Pearl, J.: *Causality: Models, Reasoning and Inference*. Cambridge University Press, Cambridge (2000)
21. Perinotti, P.: *Causal structures and the classification of higher order quantum computations* (2016)
22. Pienaar, J., Brukner, Č.: A graph-separation theorem for quantum causal models. *New J. Phys.* **17**(7), 073020 (2015)
23. Ried, K., Agnew, M., Vermeyden, L., Janzing, D., Spekkens, R.W., Resch, K.J.: A quantum advantage for inferring causal structure. *Nat. Phys.* **11**, 1745–2473 (2015)
24. Selinger, P.: A survey of graphical languages for monoidal categories. In: Coecke, B. (ed.) *New Structures for Physics*. LNP, vol. 813. Springer, Heidelberg (2011)
25. Shpitser, I., Pearl, J.: Identification of joint interventional distributions in recursive semi-Markovian causal models. In: *Proceedings of the National Conference on Artificial Intelligence*, vol. 21, p. 1219. AAAI Press/MIT Press, Menlo Park/Cambridge (1999/2006)
26. Tian, J., Pearl, J.: A general identification condition for causal effects. In: *Proceedings of the Eighteenth National Conference on Artificial Intelligence and Fourteenth Conference on Innovative Applications of Artificial Intelligence*, 28 July–1 August 2002, Edmonton, Alberta, Canada, pp. 567–573 (2002)
27. Wootters, W.K., Zurek, W.H.: A single quantum cannot be cloned. *Nature* **299**(5886), 802–803 (1982)

Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

