

Article

# Towards a Framework for Observational Causality From Time Series: When Shannon Meets Turing

David Sigtermans<sup>1\*</sup>

<sup>1</sup> ASML

\* Correspondence: david.sigtermans@asml.com

**Abstract:** We propose a novel tensor-based formalism for inferring causal structures from time series. An information theoretical analysis of *transfer entropy* (*TE*), shows that *TE* results from transmission of information over a set of communication channels. Tensors are the mathematical equivalents of these multi-channel causal channels. A multi-channel causal channel is a generalization of a *discrete memoryless channel* (*DMC*). We consider a *DMC* as a single-channel causal channel. Investigation of a system comprising three variables shows that in our formalism, bivariate analysis suffices to differentiate between direct and indirect relations. For this to be true, we have to combine the output of multi-channel causal channels with the output of single-channel causal channels. We can understand this result when we consider the role of noise. Subsequent transmission of information over noisy channels can never result in less noisy transmission overall. This implies that a *Data Processing Inequality* (*DPI*) exists for transfer entropy.

**Keywords:** information theory; transfer entropy; time-delayed mutual information; data processing inequality; time series; causal tensor

---

## 1. Introduction

### 1.1. Motivation and Significance of the Work

Exact knowledge about the causal relationships that determine the behavior of complex systems is a holy grail in the (applied) sciences and engineering. This knowledge enables us to determine potential causes of certain effects, and it allows us to predict the effect of changes in causes (for example via simulation). In other words, it allows for causal inference and causal discovery [1]. In a causal relation, the cause precedes the effect (temporal precedence), and the cause physically influences the effect [2]. A causal description is essentially different from a description via statistical associations as illustrated by the adage “*correlation does not imply causation*”. Examples of wrong, expensive, or even worse, harmful conclusions and policies based on statistical associations are part of common lore.

For a causal description, intervention is required [3]. These interventions enable us to differentiate between direct and indirect, or spurious, associations<sup>1</sup>. Because interventions are not always possible, we have to make do with observational data. A plethora of methods to infer causal structures from observational data have been developed, see for example [4–13]. What most these methods have in common is that they express relations via the causal effect, i.e., point-wise estimators that characterize “strength” of the association between a cause and an effect.

We propose a novel approach inspired by Turing machines [14]. If a human “computer” can decide, given the data, if a relation is causal, a Turing machine exists that reaches this decision in a

---

<sup>1</sup> An indirect association is an association via one or more mediators.

32 mechanical way [15]. This is not a tautology, a Turing machine encodes the underlying principles  
33 leading to the decision that a relation is causal. It closely relates our approach to Structural Causal  
34 Models [16]. Instead of using point-wise estimators like transfer entropy [8] or time-delayed mutual  
35 information [5], we use stochastic tensors, i.e. multilinear maps.

36 Under the assumption that there are no hidden causes, unmeasured common causes or  
37 confounders, our formalism can differentiate between direct and indirect associations. We show  
38 that noise has a fundamental and, from the viewpoint of detecting spurious associations, a functional  
39 role. It is as if noise acts like “soft” interventions [17]. A surprising result is that for time-delayed  
40 mutual information *bivariate analysis suffices* to differentiating between direct and indirect associations.  
41 This result contradicts the long-held belief that this is impossible (see for example, [18] and [19]). The  
42 formalism furthermore allows for a simple proof of a Data Processing Inequality [20] for transfer  
43 entropy. This DPI can identify potential indirect relations when using TE, see for example [21].

#### 44 1.2. Outline

45 The proposed formalism relies heavily on probability theory [22] and aspects and terminology  
46 used in causal inference [3]. In Section 1.3 we give a short overview of the most important ones.

47 To derive our formalism, we apply concepts from information theory [23] to transfer entropy.  
48 Transfer entropy is a measure that *can* capture causal relations, as far as encoded in the probability  
49 density functions, see for example [24,25]. In Section 2 we therefore introduce the applicable  
50 aspects of information theory, e.g., transmission of information, mutual information, communication  
51 channels and the tensor representation of communication channels. A tensor is a multilinear map that  
52 transforms an input into an output. We introduce transfer entropy in Section 3. We then show that  
53 transfer entropy is the average mutual information resulting from transmission of information of a set  
54 of communication channels. We call this set of channels *multi-channel causal channels*. Tensors are the  
55 mathematical equivalent to the set of channels. Using these tensors, we establish calculation rules in  
56 Section 4. We restrict ourselves to a system comprising three variables. We derive these results using  
57 index notation in Section 4. The result allow for a different notation that helps us to avoid a notational  
58 jungle of indices. This notation is borrowed from quantum mechanics and we introduce it in Section  
59 5. This let us incorporate the temporal relations between a source and a destination, reflecting the  
60 additivity of interaction delays. Using this notation, we discuss some relevant findings for causal  
61 inference in Section 6.

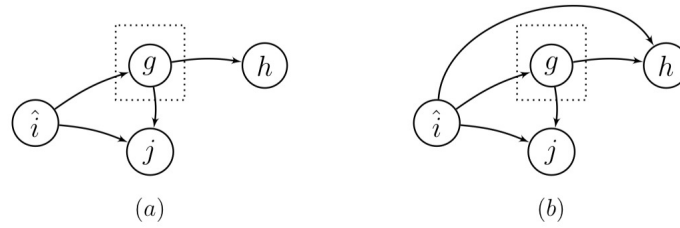
62 In Section 7 we present two experiments that illustrate that our formalism can indeed detect  
63 nonlinear relationships and an underlying structure.

64 For readability, we have moved the longer proofs or sketches of proofs to the Appendix.

#### 65 1.3. Preliminaries

66 Statistical independence is foundational to causal inference [7]. We will summarize the two  
67 most related and relevant assumptions: (1) the faithfulness assumption. (2) the Causal Markov  
68 Condition. A directed graph is said to be faithful to the underlying probability distributions if the  
69 independence relations that follow from the graph are the same independence relations that follow  
70 from the underlying probability distributions. For the chain  $X \rightarrow Y \rightarrow Z$  the faithfulness assumption  
71 implies that  $X$  and  $Z$  are independent given  $Y$ . We denote this as  $X \perp\!\!\!\perp Z|Y$ .

72 The Causal Markov Condition states that a process is independent of its non-effects given its direct  
73 causes or its parents. This is relevant in the context of time series. A straightforward interpretation of  
74 this condition is that if the set of variables blocks all (undirected) paths between two variables, these  
75 two variables are independent given the set of variables blocking all paths [3]. We illustrate the Causal  
76 Markov Condition with an example that will be used later in this article.



**Figure 1.** The graphs used in Example 1. (a) The dotted box blocks all paths between the non-parents  $\hat{i}$  and  $j$  of  $h$ , and  $g$  is the only parent of  $h$ . (b) The dotted box does not block the path between  $\hat{i}$  and  $h$ .

77 **Example 1.** Lets start with the graph depicted in Figure 1(a). According to the Causal Markov Condition,  
78  $\{\hat{i}, j\}$  and  $h$  are independent given  $g$ :  $p(\hat{i}, j, h|g) = p(\hat{i}, j|g)p(h|g)$ .

79 We now rewrite this using expressions that follow from the definition for conditional probabilities. The  
80 left-hand side is written as  $p(\hat{i}, j, h|g) = p(j, h|\hat{i}, g)p(\hat{i}|g)$ . The right-hand side can be rewritten using  $p(\hat{i}, j|g) =$   
81  $p(j|\hat{i}, g)p(\hat{i}|g)$ .

82 This finally leads the conclusion that  $p(j, \hat{i}, g, h) = p(j|\hat{i}, g)$ . This last expression also implies that  
83  $p(j|\hat{i}, g)p(g|\hat{i}, h) = p(j, g|\hat{i}, h)$ . Both expressions will be used later in this article.

84  
85 Now consider the situation depicted in Figure 1(b). According to the Causal Markov Condition  $\{\hat{i}, j\}$  and  
86  $h$  are not independent given  $g$ , i.e.,  $p(\hat{i}, j, h|g) \neq p(\hat{i}, j|g)p(h|g)$ . We can still rewrite the left-hand side and the  
87 right-hand side in the same fashion as before.

88 This finally leads the conclusion that  $p(j, \hat{i}, g, h) \neq p(j|\hat{i}, g)$ . This implies that  $p(j|\hat{i}, g)p(g|\hat{i}, h) \neq$   
89  $p(j, g|\hat{i}, h)$ .

90 In the example above we used a simplified notation for the probabilities  $p(y) := Pr\{Y = y\}$ . As  
91 illustrated in the above example we use basic aspects of probability theory like the definition of joint  
92 probabilities and the Law of Total Probability [22]. This law links a marginal probability to a joint  
93 probability, e.g.,  $\sum_g p(j, g|\hat{i}, h) = p(j|\hat{i}, h)$ .

94 Unless stated otherwise, in this article we will make use of the Einstein summation convention  
95 (with a twist). This convention simplifies equations by implying summation over indices that appear  
96 both as upper indices and as lower indices. In our definition the summation takes place the first indices  
97 and the subsequent identical lower indices  $B_j^i A_i^j := \sum_i B_j^i A_i^j$ . Our definition implies that  $B_j^i A_i^j \neq A_i^j B_j^i$ ,  
98 the order matters.

## 99 2. Information Theory

100 Shannon introduced information theory in 1948 [23]. It models association between random  
101 variables as resulting from a communication process between a sender—the source—and a receiver  
102 or the destination. A message comprises indexed realizations of random variables representing  
103 stationary ergodic processes. An input message is first encoded: we describe the message using a finite  
104 alphabet. Each random variable has its own finite alphabet. The random variable  $X$  is mapped on  
105 symbols from the alphabet  $\mathcal{X}$ , the random variable  $Y$  is mapped on symbols from  $\mathcal{Y}$ , and the random  
106 variable  $Z$  is mapped on symbols from  $\mathcal{Z}$ . Where  $\mathcal{X} = \{\chi_1, \chi_2, \dots, \chi_{|\mathcal{X}|}\}$ ,  $\mathcal{Y} = \{\psi_1, \psi_2, \dots, \psi_{|\mathcal{Y}|}\}$ , and  
107  $\mathcal{Z} = \{\zeta_1, \zeta_2, \dots, \zeta_{|\mathcal{Z}|}\}$ . The number of elements in the alphabet—the cardinality—is denoted as  $|\mathcal{X}|$ ,  
108  $|\mathcal{Y}|$ , and  $|\mathcal{Z}|$  respectively.

109 Once encoded the message is transmitted symbol by symbol. The input symbol is transformed  
110 into an output symbol. The output alphabet can have a different cardinality than the input alphabet.  
111 The transformation from input to output symbol is modeled as a Markov chain. The probability  
112 that a specific output symbol is received only depends on the alphabet symbol that was sent. The  
113 communication process transforms the input probability mass function (pmf) into the output pmf. The

114 transmitted message is decoded and made available to the receiver. In this article, we assume that no  
115 decoding takes place.

## 116 2.1. Mutual Information

If there is an association between two messages, information is said to be shared between them. The amount of information shared,

$$I(X; Y) = \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} p(x, y) \log_2 \left[ \frac{p(y|x)}{p(y)} \right], \quad (1)$$

is nonnegative and symmetric in  $X$  and  $Y$ . This so-called mutual information (MI) represents the reduction in uncertainty about the random variable  $X$  given that we have knowledge about the random variable  $Y$  (and vice versa). It is intuitively clear that, given the information content of the source data, in subsequent transmissions, the information can never increase. This is formalized in the Data Processing Inequality or DPI which states that processing of data can never increase the amount of information [20]. For the cascade  $X \rightarrow Y \rightarrow Z$  the DPI implies that, in terms of MI,

$$I(X; Z) \leq \min[I(X; Y), I(Y; Z)].$$

117 The maximum rate with which information can be transmitted between the sender and receiver is  
118 the channel capacity  $C_{XY} = \max_{p(x)} [I(X; Y)]$ . This is achieved for a so-called channel achieving input  
119 distribution.

## 120 2.2. The Communication Channel

121 In information theory, the directed graph representing a Markov chain is represented as a  
122 *communication channel*, or channel in short. The channel has an input side—the left-hand side—and an  
123 output side—the right-hand side. On the left-hand side we place all the vertices of the Markov  
124 chain with outgoing edges and on the right-hand side we place all the vertices of the Markov  
125 chain with incoming edges. The input vertices are connected to the output vertices via undirected  
126 edges. In a channel, every input alphabet symbol has its own input vertex. Likewise, every output  
127 alphabet symbol has its own output vertex. The simplest type of channel is the noisy discrete  
128 memoryless communication channel. In a memoryless channel the output— $y_t$ —only depends on  
129 the input— $x_t$ —and not on the past inputs or outputs:  $p(y_t|x_t, x_{t-1}, y_{t-1}) = p(y_t|x_t)$ . A memoryless  
130 channel embodies the Markov property. In a noisy channel the output depends on the input and  
131 another random variable representing noise. The effect of transmitting data using a DMC is described  
132 via the Law of Total Probability because

$$Pr\{Y = \psi_j\} = \sum_i Pr\{X = \chi_i\} Pr\{Y = \psi_j | X = \chi_i\}, \quad (2)$$

133 with  $Pr\{Y = \psi_j\}$  the  $j^{\text{th}}$  element of the probability mass function  $p(y)$ , and  $Pr\{X = \chi_i\}$  the  $i^{\text{th}}$   
134 element of the pmf  $p(x)$ . The transmission of data over a DMC transforms the pmf of the input into the  
135 pmf of the output via a linear transformation. The probability transition matrix  $Pr\{Y = \psi_j | X = \chi_i\}$   
136 fully characterizes the DMC [20]. Assuming a fixed (e.g. lexicographic) order of the alphabet elements,  
137 we can introduce an index notation for the pmfs, e.g.  $p^j := Pr\{Y = \psi_j\}$  and  $p^i := Pr\{X = \chi_i\}$ . In this  
138 article we associate every index with a specific random variable. In Table 1 an overview is given.

**Table 1.** Overview of Indices Used.

Process	Variable	Alphabet element	Index (input)	Index (past)	Index (output)
X	x	$\chi$	$\hat{i}$	f	i
Y	y	$\psi$	$\hat{j}$	g	j
Z	z	$\zeta$	$\hat{k}$	h	k

### 139 2.3. Tensor Representation of the Communication Channel

140 One of the many virtues of information theory is that it enables the use of linear algebra. Because  
 141 we do not want to get overwhelmed by increasingly complex probabilistic equations, we use index  
 142 notation and the Einstein summation convention (with a minor twist). Equation (2) can now be  
 143 written as  $p^j = p^i A_i^j$ . The covariant indices indicate the variables we condition on. The row stochastic  
 144 probability transition matrix elements represent the elements of the probability transition tensor A  
 145 [26]. Using the standard notation instead of the Einstein summation convention, we can rewrite MI as  
 146  $I(X; Y) = \sum_{i,j} p^{ij} \log_2 [A_i^j / p^j]$ .

147 Mutual information solely depends on the elements of the tensor and the input pmf. This is  
 148 problematic in case MI or MI derived measures are used to infer the underlying structure if we assume  
 149 that the structure is independent of the input. We can illustrate this by assuming that the probability  
 150 transition tensor equals the Kronecker delta

$$\delta_i^j = \begin{cases} 1, & \text{if } i = j, \\ 0, & \text{if } i \neq j. \end{cases}$$

151 **Example 2.** If  $A_i^j = \delta_i^j$ , the symbol received is identical to the symbol sent, the channel transmits data perfectly.  
 152 In this case MI reduces to  $I(X; Y) = \sum_i p^i \log_2 [1/p^i]$ . Now set the probability of one of the alphabet elements  
 153 to  $1 - \epsilon$ . This implies that all other symbol probabilities are equal to or smaller than  $\epsilon$ . Taking the limit  $\epsilon \rightarrow 0$   
 154 results in a mutual information  $\rightarrow 0$ . Although there might be a noiseless channel representing the association  
 155 between the random variables X and Y, MI could be arbitrarily small.

156 This leads us to the following proposition for inferring structures using MI-based measures:

157 **Proposition 1.** In case MI or MI related measures are used to infer the **structure** for a system, we should use  
 158 the probability transition tensors or measures based on elements of probability transition tensor.

159 The earlier mentioned channel capacity is such a measure. It only depends on the elements of  
 160 the probability transition tensor [27], e.g.  $C_{XY} := \Gamma(A)$ . Because the channel capacity is the maximal  
 161 achievable mutual information for a specific channel, the earlier mentioned DPI also applies to the  
 162 channel capacity. The proof is straightforward and therefore omitted.

163 **Corollary 1** (DPI for Channel Capacity). For the chain  $X \rightarrow Y \rightarrow Z$  the DPI immediately implies that  
 164  $\Gamma(C) \leq \min[\Gamma(A), \Gamma(B)]$ , with A representing the tensor of the transmission  $X \rightarrow Y$ , B:  $Y \rightarrow Z$ , and C:  $X \rightarrow Z$ .

165 In this short and incomplete introduction to information theory, no assumptions—other than  
 166 stationarity, ergodicity and Markov property—were made about the underlying mechanisms leading  
 167 to the association between random variables. We can therefore apply information theory to all cases  
 168 where observational data are available.

### 169 3. Transfer Entropy

Schreiber introduced transfer entropy in 2000 [8]. Like MI it is non-parametric, but unlike MI  
 it is an essentially asymmetric measure and it enables the differentiation between a source and a

destination. It is an information theoretical implementation of Wiener's principle of Causality [28]: a cause combined with the past of the effect predicts the effect better than that the effect predicts itself. In contrast to Granger causality [18], transfer entropy can capture nonlinear relationships. We use a slightly modified version which was shown to comply fully with Wiener's principle of Causality by Wibral et al. They proved that this modified TE is maximal for the real interaction delay [29]. We assume that  $Y$  is a Markov process of order  $\ell \geq 1$ . This implies that the future  $y_t$  also depends on its past  $\mathbf{y}^- = (y_{t-1}, \dots, y_{t-\ell})$ . The destination also depends on the source data  $X$ . With  $\tau$  the finite interaction delay, we assume that for the input symbol  $\mathbf{x}^- = (x_{t-\tau}, \dots, x_{t-\tau-m})$ , with  $m \geq 0$ . The alphabet for the past of  $Y$  is  $\mathcal{Y}^\ell$ . The alphabet for the input is  $\mathcal{X}^m$ .

$$TE_{X \rightarrow Y} = \sum_{\substack{\mathbf{x}^- \in \mathcal{X}^m, y \in \mathcal{Y} \\ \mathbf{y}^- \in \mathcal{Y}^\ell}} p(\mathbf{x}^-, y, \mathbf{y}^-) \log_2 \left[ \frac{p(y|\mathbf{x}^-, \mathbf{y}^-)}{p(y|\mathbf{y}^-)} \right] \quad (3)$$

170

171 To differentiate a source from a destination, we have to assess two hypotheses: (1)  $X$  is the source  
172 and  $Y$  is the destination, and (2)  $Y$  is the source and  $X$  is the destination. Per case the interaction  
173 delay that maximizes the respective TE is determined. If the resulting transfer entropy equals 0, we  
174 assume that there is no relation. If the TE values are larger than 0, there are in practice two possibilities:  
175 (1) the optimal interaction delays are equal: we assume that the hypothesis resulting in the largest TE  
176 is valid. (2) The optimal interaction delays are different: both hypotheses are valid so we have detected  
177 a cycle. Without loss of generality, we assume in this article that there are no cycles.

178 Transfer entropy is a conditional mutual information [8]. Therefore, it can be associated with  
179 communication channels. We start with conditioning the MI from Equation (1) on the event  $\mathbf{y}^- = \psi_g^-$   
180 resulting in

$$I(X; Y | \psi_g^-) = \sum_{\substack{\mathbf{x}^- \in \mathcal{X}^m \\ y \in \mathcal{Y}}} p(\mathbf{x}^-, y | \psi_g^-) \log_2 \left[ \frac{p(y|\mathbf{x}^-, \psi_g^-)}{p(y|\psi_g^-)} \right].$$

Because  $\mathbf{x}^-$  and  $\mathbf{y}^-$  are the only parents of the output  $y$ , it follows from the Causal Markov Condition that the associated channel is memoryless. This mutual information quantifies the amount of information that is transmitted over the  $g^{\text{th}}$  sub-channel. The transfer entropy from Equation (3) can now be expressed as

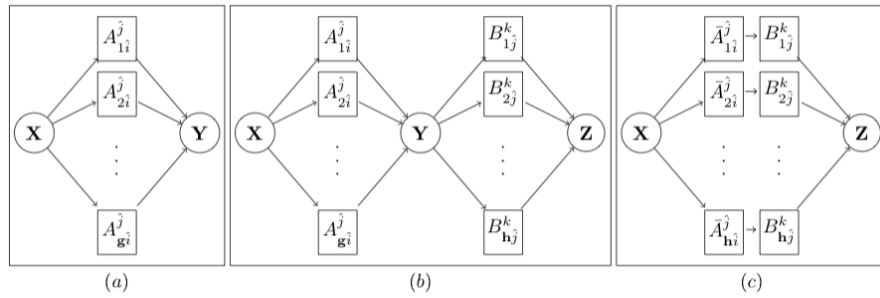
$$TE_{X \rightarrow Y} = \sum_{\psi_g^- \in \mathcal{Y}^\ell} p(\psi_g^-) I(X; Y | \psi_g^-). \quad (4)$$

181 We now show that this expression results from transmission of information over a set of  
182 communication channels.

### 183 3.1. The Causal Channel

184 An inverse multiplexer comprises a demultiplexer and a multiplexer in series. A demultiplexer  
185 separates an input data stream into multiple output data streams. We call these different streams  
186 sub-channels. A multiplexer combines or multiplexes all input data streams into a single output data  
187 stream [30].

188 **Definition 1** (Causal Channel). *A causal channel is an inverse multiplexer in which the demultiplexer selects*  
189 *the sub-channel over which the data are sent based on the past of the output data. Each sub-channel consists of a*  
190 *DMC. The input symbol is fed to a specific input vertex of the chosen discrete memoryless channel. The DMC*  
191 *transforms the input in a probabilistic fashion into an output symbol. The multiplexer combines the outputted*  
192 *symbols into the output message. See Figure 2(a).*



**Figure 2.** (a) Causal channel. (b) Two causal channels in series representing the communication model related to transfer entropy for the cascade  $X \rightarrow Y \rightarrow Z$ . (c) The equivalent causal channel for two causal channels in series.

193 This definition forms the basis for the theorem that is central to this article.

194 **Theorem 1.** *Transfer entropy is the average conditional mutual information of transmission over a causal*  
 195 *channel.*

196 **Proof.** The relative frequency with which the  $g^{th}$  sub-channel is chosen equals  $p(\psi_g^-)$ . Each  
 197 sub-channel is a DMC, so the mutual information of the  $g^{th}$  sub-channel equals  $I(X; Y | \psi_g^-)$ .  
 198 The weighted average of the mutual information over all the sub-channels is equal to  
 199  $\sum_{\psi_g^- \in \mathcal{Y}^\ell} p(\psi_g^-) I(X; Y | \psi_g^-)$ , which is the definition of TE in Equation (4).  $\square$

200 Because a DMC is a causal channel with only one sub-channel, we call a DMC a *single-channel*  
 201 *causal channel*.

### 202 3.2. Tensor Representation of a Causal Channel

203 Because every sub-channel of the causal channel represents a DMC, a causal channel can be  
 204 represented by a probability transition tensor. We will call this tensor a *causal tensor*. For the relation  
 205  $X \rightarrow Y$  we get the following equation for the  $g^{th}$  sub-channel

$$p_g^j = p_g^i A_{gi}^j. \quad (5)$$

The elements of the tensor A are given by  $A_{gi}^j = p(\psi_j | \chi_i^-, \psi_g^-)$ . We can now rewrite TE as

$$TE_{X \rightarrow Y} = \sum_{g, \hat{i}, j} p^{i j g} \log_2 \left[ \frac{A_{gi}^j}{p_g^j} \right]. \quad (6)$$

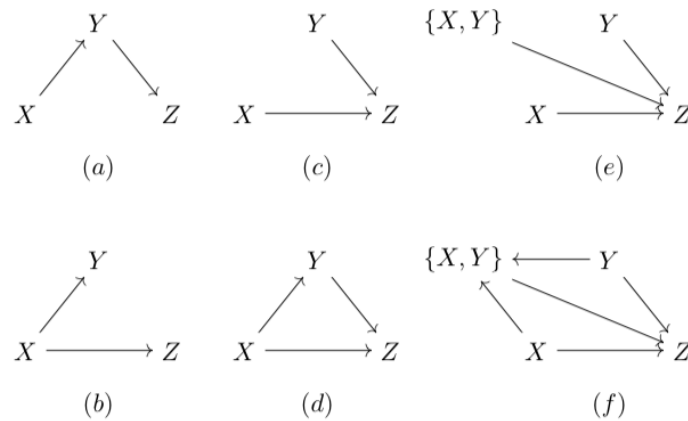
206 In a similar fashion as MI, we can show that TE can be made arbitrarily close to 0 while the  
 207 causal tensor itself represents a noiseless transmission. It is therefore not an optimal measure to infer  
 208 structures. Again we would prefer to use the tensors themselves or measures based on these tensors  
 209 like the channel capacity. The determination of the channel capacity for a causal channel is not in the  
 210 scope of this article. We assume however that it is possible to determine the channel capacity.

211 As stated in the introduction, the approach in the article was inspired by Turing machines. The  
 212 causal tensor is a realization of the transition function of a Turing machine that encodes causality in as  
 213 far as the causality is encoded in the pmfs. To warrant the use of the adjective “causal” however, we  
 214 have to show that within the framework of causal tensors; we can differentiate between direct and  
 215 indirect associations. That this seems possible can be intuited when considering the chain  $X \rightarrow Y \rightarrow Z$   
 216 (see Figure 3a). The relation  $X \rightarrow Z$  is a resultant of the other relations, i.e., an indirect association.

217 Within the framework of causal tensors, we would expect that we can express this indirect association  
 218 in terms of the tensors of the other relations.

#### 219 4. Calculation Rules for Causal Tensors

220 In this section we derive the calculation rules for causal tensors. Operations performed on these  
 221 tensors should result in either scalars, stochastic vectors or stochastic tensors. To prove that we end up  
 222 with stochastic tensors or vectors, we will use the earlier introduced index notation, the law of total  
 223 probability and the Causal Markov Condition. We derive the calculation rules by investigating the  
 224 four elementary structures depicted in Figures 3(a), 3(b), 3(c) and 3(d).



**Figure 3.** The basic structures directed graph structures: (a) the chain, (b) the fork, (c) the v-structure, and (d) the directed triangle. The graphs (e) and (f) reflect the calculation rules for the causal tensors for the v-structure and directed triangle respectively.

225 In Section 5 we will propose another notation which is simpler, but it relies on the results in this  
 226 section.

##### 227 4.1. The Chain Structure

228 First assume the chain  $X \rightarrow Y \rightarrow Z$  is the ground truth. Additional to Equation (5),  $p_g^j = p_g^i A_{g_i}^j$ ,  
 229 there are two other causal channels represented by tensors:  $B: Y \rightarrow Z$ , and  $C: X \rightarrow Z$ . Because it is a  
 230 straightforward exercise in which we again make use of the law of total probability, we leave it to the  
 231 reader to confirm that

$$p_h^k = p_h^{\hat{i}} B_{h\hat{i}}^k \quad (7a)$$

$$p_h^k = p_h^{\hat{i}'} C_{h\hat{i}'}^k. \quad (7b)$$

232 In principle the index  $\hat{i}'$  in Equation (7b) is the index representing a different input vector than  
 233 the index  $\hat{i}$  in Equation (5); although they both refer to the random variable  $X$ . This is because  $\hat{i}'$  is  
 234 related to the source  $\mathbf{x}' \in \mathcal{X}^{m'}$  of  $Z$  and  $\hat{i}$  is the index related to the source  $\mathbf{x} \in \mathcal{X}^m$  of  $Y$ . The Markov  
 235 property however immediately implies that we can use one and the same index in both cases if we  
 236 select the source vector with the largest cardinality. Without loss of generality, we use  $\hat{i}$ .

237 We can express the tensor of the indirect relation in terms of the tensors of the direct relations.



238 **Theorem 2** (Product Rule for a Chain). *Let A and B be the causal tensors of two causal channels in series.*  
 239 *Let the tensor C represent the resulting indirect causal channel that must be measured in a bivariate approach.*  
 240 *The tensor elements of C are given by*

$$C_{hi}^k = p_{hi}^g A_{gi}^j B_{hj}^k. \quad (8)$$

241 *If the ground truth is a directed triangle, the product rule for a chain is invalid, i.e.,  $C_{hi}^k \neq p_{hi}^g A_{gi}^j B_{hj}^k$ .*

242 The proof is given in Appendix A.1. The term  $p_{hi}^g A_{gi}^j$  is very interesting. In Appendix A.2 we  
 243 proof that it is a stochastic tensor.

244 **Lemma 1.** *For a chain the product  $p_{hi}^g A_{gi}^j$  is a stochastic tensor  $\bar{A}_{hi}^j$ .*

245 The causal tensor  $p_{hi}^g A_{gi}^j$  can be interpreted as the weighted average of the causal tensor A, given  
 246 the  $h^{th}$  sub-channel of the final causal channel and the input  $\hat{i}$ . We can now rewrite Equation (8) as

$$C_{hi}^k = \bar{A}_{hi}^j B_{hj}^k. \quad (9)$$

247 If both A and B represent discrete memoryless channels we get the simpler, well known, product  
 248 rule for a chain of DMC's.

249 **Corollary 2** (Product Rule for a Chain of DMC's). *Let A and B be the causal tensors of two DMC's in*  
 250 *series and let the tensor C represent the resulting, indirect, causal channel that must be measured in a bivariate*  
 251 *approach. The tensor elements of C are given by*

$$C_{\hat{i}}^k = A_{\hat{i}}^j B_{\hat{j}}^k. \quad (10)$$

252 The proof follows directly from the definition of a DMC in terms of a causal channel: a discrete  
 253 memoryless channel is a causal channel comprising only one sub-channel, i.e., it is a single-channel  
 254 causal channel. Combined with Lemma 1, this corollary leads to a very specific interpretation of  
 255 Equation (8). According to Corollary 2, Equation 9 can be interpreted as representing two DMC's in  
 256 series for the  $h^{th}$  sub-channel. This means we have an alternative structure for two causal channels in  
 257 series as depicted in Figure 2(c).

258 Because the Data Processing Inequality applies to a cascade of discrete memoryless channels,  
 259 the alternative structure suggests that there is a DPI for transfer entropy. In Section 6 we show that  
 260 this is indeed the case. If one so wishes we could check if the measured TE for the potential spurious  
 261 association equals the expected TE. For the chain  $X \rightarrow Y \rightarrow Z$  the expected  $TE_{X \rightarrow Z}$  is given by

$$TE_{X \rightarrow Z} = \sum_{\hat{i}, h, k} p^{\hat{i}hk} \log_2 \left[ \frac{\sum_j \bar{A}_{hi}^j B_{hj}^k}{p_h^k} \right].$$

#### 262 4.2. The Fork Structure

263 In this section, we show that a fork can be interpreted as a chain. The product rule for a chain is  
 264 therefore also applicable to a fork. Assume that the fork is the ground truth (Figure 3(d)). Again we  
 265 want to express the indirect association represented by B in terms of the other causal tensors. First, we  
 266 notice that the input distribution can be reconstructed from the output distribution.

267 **Definition 2** (Reconstruction Operator). *The  $\ddagger$ -operator, or reconstruction operator, reconstructs the source*  
 268 *distribution, conditioned of the past of the destination, from the destination distribution, conditioned of the past*  
 269 *of the destination:*

$$p_g^{\hat{i}} = p_g^j A_{gj}^{\ddagger \hat{i}}, \quad (11)$$

270 *with  $A_{gj}^{\ddagger \hat{i}} = p_{gj}^{\hat{i}}$ . The  $\ddagger$ -operation changes the sign of the interaction delay of the original relation.*

271 This implies that the directed graph  $X \rightarrow Y$ , is equivalent to the graph  $X \leftarrow \ddagger Y$ . Because it is  
 272 straightforward using Equation (6) we leave it to the reader to prove the following corollary.

273 **Corollary 3.** *From an information theory point of view, a relation and its reconstructed relation are equivalent,*  
 274 *i.e.,  $TE_{X \rightarrow Y} = TE_{X \leftarrow \ddagger Y}$*

275 From this corollary immediately follows that a fork is equivalent to a chain.

276 **Theorem 3** (Fork-Chain Equivalence). *The fork  $X \rightarrow Y + X \rightarrow Z$  is equivalent to the chain  $Y \ddagger \rightarrow X \rightarrow Z$  and*  
 277 *to the chain  $Y \leftarrow X \leftarrow \ddagger Z$ .*

278 The indirect association represented by B in terms of the other two tensors of the chain follows  
 279 directly from the product rule for a chain (Theorem 2).

$$B_{h\hat{j}}^k = \bar{A}_{h\hat{j}}^{\ddagger \hat{i}} C_{h\hat{i}}^k, \text{ with } \bar{A}_{h\hat{j}}^{\ddagger \hat{i}} := p_{h\hat{j}}^g A_{g\hat{j}}^{\ddagger \hat{i}}, \quad (12a)$$

$$B_{g\hat{k}}^j = \bar{C}_{g\hat{k}}^{\ddagger \hat{i}} A_{g\hat{i}}^j, \text{ with } \bar{C}_{g\hat{k}}^{\ddagger \hat{i}} := p_{g\hat{k}}^h C_{h\hat{k}}^{\ddagger \hat{i}}. \quad (12b)$$

280 Equation (12a) applies in case the equivalent chain is  $Y \rightarrow \ddagger X \rightarrow Z$ . If the equivalent chain is  
 281  $Y \leftarrow X \leftarrow \ddagger Z$ , Equation (12b) is applicable. Due to the way we determine the interaction delay, the  
 282  $\ddagger$ -operation induces a sign change for the interaction delay. E.g., if  $\tau_{xy}$  represents the interaction delay  
 283 for the relation  $X \rightarrow Y$ , then  $-\tau_{xy}$  represents the interaction delay for the relation  $Y \ddagger \rightarrow X$ .

284 Matrices, and therefore tensors, do not commute:  $\bar{A}_{h\hat{j}}^{\ddagger \hat{i}} \bar{A}_{h\hat{i}}^{\ddagger \hat{j}} \neq \bar{A}_{h\hat{i}}^{\ddagger \hat{j}} \bar{A}_{h\hat{j}}^{\ddagger \hat{i}}$ . This leads to the conclusion  
 285 that a chain and a fork are in principle distinguishable. The reader can verify this by combining  
 286 Equation (9) and Equation (12a). Combining this with Theorem 2 leads to the conclusion that we can  
 287 differentiate between direct and indirect relations. The conditions under which it is not possible will  
 288 be derived later.

289 **Theorem 4.** *Using causal tensors we can differentiate between a chain, a fork and a directed triangle. If and*  
 290 *only if the chain is the ground truth  $B_{h\hat{j}}^k \neq \bar{A}_{h\hat{j}}^{\ddagger \hat{i}} C_{h\hat{i}}^k$ . If and only the fork is the ground truth  $C_{h\hat{i}}^k \neq \bar{A}_{h\hat{i}}^{\ddagger \hat{j}} B_{h\hat{j}}^k$ . If the*  
 291 *structure is neither a chain, nor a fork, it is a directed triangle.*

292 In case of a single-channel causal channels, we are allowed to ignore the index indicating the  
 293 sub-channels. The equations for the chain, the fork, and consequently the equations in Theorem 4 are  
 294 all bivariate.

295 **Corollary 4.** *If we use single-channel causal tensors, bivariate time-delayed mutual information measurements*  
 296 *can differentiate between a fork, a chain, and a directed triangle.*

297 This result contradicts the current point of view [12,19]. We illustrate this with an example.  
 298 Because we use single-channel causal channels, we ignore the index  $h$  in the equations.

299 **Example 3.** Let the chain  $X \rightarrow Y \rightarrow Z$  be the ground truth. With  $A = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ 1 & 0 \end{pmatrix}$  and  $B = \begin{pmatrix} \frac{1}{3} & \frac{2}{3} \\ 0 & 1 \end{pmatrix}$ , the indirect  
 300 association is represented by the causal tensor  $C_{\hat{i}}^k = A_{\hat{i}}^j B_{\hat{j}}^k \Rightarrow C = \begin{pmatrix} \frac{1}{6} & \frac{5}{6} \\ \frac{1}{3} & \frac{2}{3} \end{pmatrix}$ . Assume that  $p(x) = (\frac{2}{5}, \frac{3}{5})$ . The  
 301 pmf for  $p(y)$  equals  $p(y) = p(x)A$ . From this follows that  $p(y) = (\frac{4}{5}, \frac{1}{5})$ . Using the relation  $p(x) = p(y)A^\dagger$ , the  
 302 reader can verify that  $A^\dagger = \begin{pmatrix} \frac{1}{4} & \frac{3}{4} \\ 1 & 0 \end{pmatrix}$ . Because  $\bar{A}_{\hat{h}\hat{i}}^\dagger C_{\hat{h}\hat{i}}^k = \begin{pmatrix} \frac{7}{24} & \frac{17}{24} \\ \frac{1}{6} & \frac{5}{6} \end{pmatrix}$ . From this follows that  $B_{\hat{h}\hat{j}}^k \neq \bar{A}_{\hat{h}\hat{j}}^\dagger C_{\hat{h}\hat{i}}^k$ . The  
 303 structure is that of a chain, and not that of a fork.

#### 304 4.3. The V-Structure and the Directed Triangle

In a bivariate measurement, we will always be able to determine the ground truth correctly in the case of the v-structure depicted in Figure 3(c). However, investigating structures with a collider, the v-structure and the more general directed triangle, will result in the important concept of *interaction*. So, let us assume that the ground truth is the directed triangle. We now have to introduce the multivariate relation  $D: \{X, Y\} \rightarrow Z$ . This relation leads to the additional linear transformation

$$p_h^k = p_h^{\hat{i}\hat{j}} D_{\hat{h}\hat{i}\hat{j}}^k.$$

305 We call the tensor D the *interaction tensor*. The tensors B and C can be expressed in terms of the  
 306 tensor D.

307 **Lemma 2** (Causal Tensor Contraction). *In the case of a directed triangle, we can express the causal tensors in*  
 308 *terms of the interaction tensor:*

$$B_{\hat{h}\hat{j}}^k = \bar{A}_{\hat{h}\hat{j}}^\dagger D_{\hat{h}\hat{i}\hat{j}}^k, \quad (13a)$$

$$C_{\hat{h}\hat{i}}^k = \bar{A}_{\hat{h}\hat{i}}^\dagger D_{\hat{h}\hat{i}\hat{j}}^k. \quad (13b)$$

309

310 For the proof, we use the fact that the elements of a causal tensor are conditional probabilities.  
 311 Due to the fork-chain equivalence, Appendix A.3 only contains the proof for the chain.

312 From Equation (13) it follows that B and C are the result of a cascade involving  $A^\dagger$  and D for  
 313 B, and A and D for C. The graphs represented by Figures 3(e) and 3(f) support the tensor relations,  
 314  $X \rightarrow \{X, Y\} \rightarrow Z$  is equivalent to the cascade of the inverse multiplexers represented by A and D  
 315 resulting in C. Figures 3(c) and 3(d) however do not support the calculation rules for causal tensors.

316 **Proposition 2.** *If a complex system contains v-structures, the causal graph must be represented by a directed*  
 317 *hypergraph [31]. In a hypergraph, an edge connects any number of vertices. The interaction tensor corresponds*  
 318 *to a so-called hyperedge.*

319 The interaction tensor describes the interaction of inputs at the v-structure. An indirect relation  
 320 does not interact.

321 **Theorem 5.** *The interaction tensor only depends on the direct causes, not on indirect causes. So, if and only if*  
 322 *the chain is the ground truth  $D_{\hat{h}\hat{i}\hat{j}}^k = B_{\hat{h}\hat{j}}^k$ . If and only if the fork is the ground truth  $D_{\hat{h}\hat{i}\hat{j}}^k = C_{\hat{h}\hat{i}}^k$ .*

323 **Sketch of Proof.** Let the ground truth be the chain. In that case,  $X \perp\!\!\!\perp Z|Y$  and X is a non-effect of Z.  
 324 The index  $\hat{i}$  is associated with X, the index  $\hat{j}$  is associated with Y and the indices  $h$  and  $k$  are associated  
 325 with Z. The Causal Markov Condition leads to  $p_{\hat{h}\hat{j}}^{\hat{i}k} = p_{\hat{h}\hat{j}}^k p_{\hat{h}\hat{j}}^{\hat{i}}$   $\Leftrightarrow p_{\hat{h}\hat{i}\hat{j}}^k = p_{\hat{h}\hat{j}}^k$ .  $\square$

326 In this index ridden section, we have shown that we can express indirect relations in terms of  
 327 the direct relations; the resulting tensors are stochastic tensors. We can now introduce a notation that  
 328 simplifies our expressions.

## 329 5. A New Notation

330 The input for a causal tensor is a probability vector. This vector represents the *probabilistic state* of  
 331 the input. Let us assume that the input and output alphabets each consists of five alphabet elements.  
 332 Because we place the operand before the operator, the probability vectors are row vectors. An input  
 333 event  $x$  could be represented by the probability vector  $(0, 0, 1, 0, 0)$ , i.e., the input equals the 3rd  
 334 alphabet element. After the transmission over the causal channel, the probabilistic *state* of the output  
 335 element could, for example, be  $(0, 0.3, 0, 0.2, 0.5)$ , a probabilistic mix of the output alphabet elements.

336 Apart from the probabilistic state, the states of the input and output events should also  
 337 accommodate the *temporal* relationship between the input and output, indicated by the interaction  
 338 delay. We can unify both state characteristics by introducing complex state vectors and by adopting  
 339 the *braket* or Dirac notation used in quantum mechanics [32].

### 340 5.1. Braket Notation and Tensor Operations

341 The complex state vector, or the *ket* is represented by the notation  $|x\rangle$ . From the viewpoint of  
 342 index notation, the order of the operand and operator is irrelevant, so we chose to represent the  
 343 probability component as a column vector, i.e., the transpose of the probability vector. In the case of  
 344 our example this would be  $(0, 0, 1, 0, 0)^T$ . The temporal dependencies are taken care of via a complex  
 345 number, e.g.,  $e^{i\omega t_x}$ . The  $i$  in  $e^{i\omega t_x}$  is not an index, but the imaginary number  $i$ . The  $\omega$  represents a  
 346 normalizing constant and  $t_x$  a time constant time related to the random variable, in this case,  $X$ . Each  
 347 random variable has its own constant time associated with it. The ket in our example is represented by  
 348  $|x\rangle = (0, 0, 1, 0, 0)^T e^{i\omega t_x}$ .

349 We assume that the interaction delay is induced by the causal channel. We therefore associate it  
 350 with the causal tensor. We represented it by  $e^{i\omega \tau_{xy}}$ , where  $\tau_{xy}$  is the interaction delay between  $X$  and  $Y$ .  
 351 The *complex* causal tensor is defined as  $\mathcal{A} := A^T e^{i\omega \tau_{xy}}$ . We can now rewrite Equation (5) as  $|y\rangle = \mathcal{A} |x\rangle$ .  
 352 In the case of our example, the ket for the output is given by  $|y\rangle = (0, 0.3, 0, 0.2, 0.5)^T e^{i\omega(t_x + \tau_{xy})}$ . To be  
 353 able to simplify the product rule for a chain given by Equation (8), we define the cascading operator  
 354 using the chain  $X \rightarrow Y \rightarrow Z$ .

355 **Definition 3.** The causal tensor cascading operator  $\odot$  applied to a cascade of two causal tensors,  $\mathcal{A}$  and  $\mathcal{B}$ , is  
 356 defined as

$$\mathcal{B} \odot \mathcal{A} |x\rangle = |z\rangle.$$

357

358 Equation (9) can now be written as

$$\mathcal{C} = \mathcal{B} \odot \mathcal{A}. \quad (14)$$

359 Our definition implies that the interaction delay of a cascade of causal channels is the sum of all  
 360 the interaction delays between the subsequent pairs making up the cascade. This is not an uncommon  
 361 assumption [33,34]. The additivity also applies to “dagged” relations.

**Definition 4.**

$$\mathcal{A}^\dagger := (\mathcal{A}^\dagger)^T e^{-i\omega \tau_{xy}}.$$

362 Apart from the fact that index ridden equations can be simplified, and independent of our  
 363 historical attachment as the framework resulted from the intuition that a notation like this should be  
 364 possible, it is a subject of future research if this novel notation provides additional new insights.

### 365 5.2. Transfer Entropy and Mutual Information

366 Both TE and MI can be expressed in terms of the complex tensors. First, we introduce a simplified  
 367 notation for TE We write these measures as a function of pmfs, indicated by  $(\cdot)$  and the respective  
 368 tensor.

#### Corollary 5.

$$TE_{X \rightarrow Y} = \|TE(\mathcal{A}, \cdot)\|,$$

369 with  $\|\cdot\|$  the Euclidean norm.

370 The proof is straightforward when using Equation (6) and we therefore omit it. From now on we  
 371 always assume that for measures like TE and MI we have to take the Euclidian norm. This allows  
 372 us to write  $TE_{X \rightarrow Y} = TE(\mathcal{A}, \cdot)$ . We use the subscripts, e.g.,  $h$ , to indicate the slice representing the  
 373 sub-channel in a multi-channel causal channel, e.g.,  $I(\mathcal{A}_h, \cdot)$  represents the MI for the  $h^{th}$  sub-channel.

## 374 6. Inferring Structures With Causal Tensors

375 In this section, we discuss some non-trivial implications when using causal tensors to infer  
 376 the causal structure from time-series data. First, we will show that a Data Processing Inequality  
 377 for transfer entropy exists. Because we did not make any assumption about the cardinality of the  
 378 alphabets used, this DPI is also valid for time-discrete, continuous data. We then prove that we can  
 379 differentiate between a fork, a chain and a directed triangle as long as the data are noisy, but not  
 380 “perfectly noisy”—we will define later this in this article. Finally we will establish a theorem that  
 381 enables us to always use bivariate analysis within a system comprising three variables.

### 382 6.1. Data Processing Inequality for TE

383 The DPI for TE gives a sufficient condition to assess if a relation is a proper direct relation. It gives  
 384 a necessary condition to detect potential indirect relations.

385 **Theorem 6** (DPI for a Chain). *For the chain  $X \rightarrow Y \rightarrow Z$  the following inequality holds*

$$TE_{X \rightarrow Z} \leq \min [TE_{X \rightarrow Y}, TE_{Y \rightarrow Z}]. \quad (15)$$

386 *Because the fork has equivalent chains, the DPI also applies to a fork.*

387 For readability, we have moved the proof to Appendix A.4. Under the condition that the  
 388 embedding of the cause and the effect vectors is sufficiently large, the DPI can identify potential  
 389 indirect relations. Because we made no assumptions about the cardinality of the (finite) alphabets, the  
 390 DPI is also valid for finite, very large alphabets.

### 391 6.2. Differentiating Between Direct and Indirect Associations With Causal Tensors

392 We have shown earlier that a fork, a chain, and a directed triangle are distinguishable. We now  
 393 investigate in more detail under what conditions this is not possible.

394 **Definition 5** (Perfect Noisy Relation). *If and only if all causal tensor elements are equal, the relation is a*  
 395 *perfect noisy relation. The related causal tensor is called the perfect noisy causal tensor.*

396 The behavior of a perfect noisy causal tensor is straightforward and therefore left to the reader  
 397 to confirm: (1) any input pmf is transformed into a uniform probability distribution, (2) the channel  
 398 capacity = 0. The opposite of the perfect noisy causal tensor is the noiseless causal tensor.

399 **Definition 6** (Noiseless Causal Tensor). *The elements of a noiseless causal tensor satisfy  $\forall_{hi\hat{j}} A_{hi}^{\hat{j}} \in \{0, 1\} \cup$*   
 400  *$\forall_h : \sum_{\hat{i}} A_{hi}^{\hat{j}} = 1$  and  $\sum_{\hat{j}} A_{hi}^{\hat{j}} = 1$ .*

401 The reader can verify by using Equation (6) that for any input pmf,  $TE = \log_2 \left[ \sum_{\hat{j}} 1 \right]$ . Because  
 402 the channel capacity of a noiseless channel only depends on the number of alphabet elements,  $C_{XY} =$   
 403  $\min \left[ \log_2(|\mathcal{X}^m|), \log_2(|\mathcal{Y}^\ell|) \right]$  [20], our definition is indeed a noiseless causal channel. An immediate  
 404 consequence of the definition of a noiseless tensor is that the cardinality of the input pmf equals the  
 405 cardinality of the output pmf.

406 That we can differentiate between direct and indirect relations is related to noise. We proved the  
 407 following theorem in Appendix A.5.

408 **Theorem 7.** *We cannot differentiate between direct and indirect relations if: (1) all relations are perfectly*  
 409 *noiseless, or (2) the relations are (almost) perfectly noisy.*

### 410 6.3. Bivariate Analysis With Causal tensors

411 We now show that within our simple system of three variables, bivariate analysis suffices. We  
 412 first need to determine the causal tensors representing multi-channel causal channels, after which we  
 413 determine the causal tensors representing single-channel causal channels.

414 **Theorem 8.** *If in a system comprising three variables, the ground truth is a chain, then the product rule for a*  
 415 *chain is applicable to both the multi-channel causal channels and the single-channel causal channels:*

$$416 \quad C_{hi}^k = \bar{A}_{hi}^{\hat{j}} B_{hj}^k \Leftrightarrow C_i^k = A_i^{\hat{j}} B_j^k. \quad (16)$$

The proof can be found in Section A.6.

417 From the proof it immediately follows that this theorem is only valid when we use the same  
 418 embedding and the same interaction delays for both the single-channel causal channels as for the  
 419 multi-channel causal channels.

### 420 6.4. Causal Inference Steps

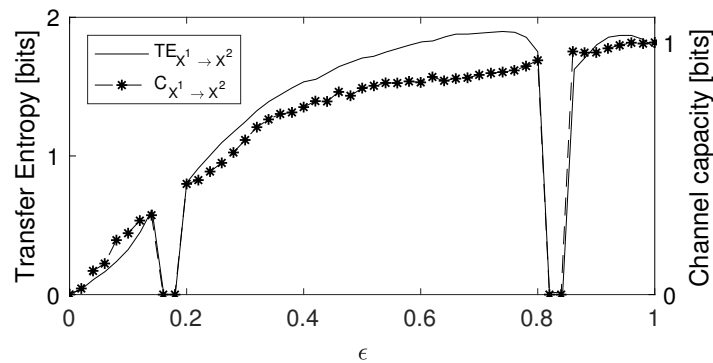
421 To complete the causal tensor framework as discussed so far, a summary of the (implicitly)  
 422 proposed steps are given. We assume that: (1) the data are time equidistant, (2)  $\ell$  and  $m$  are determined  
 423 correctly, and (3) the data are ergodic and stationary.

- 424 1. Encode the data into a finite alphabet.
- 425 2. Determine the (bivariate) multi-channel causal tensors for a range of interaction delays.
- 426 3. Determine the optimal interaction delay.
- 427 4. Determine per relation the direction of causation.
- 428 5. Identify the potential indirect relations using the DPI and the additivity of interaction delays.
- 429 6. Determine the single-channel causal tensors for the potential indirect relations.
- 430 7. Use the product rule to determine if the indirect relations are indeed indirect.
- 431 8. If the network is used for simulation, determine the interaction tensors for all v-structures, i.e.,  
 432 determine the hypergraph.

## 433 7. Experiments

434 We finalize this article with two experiments to illustrate that nonlinear behavior is indeed  
 435 captured with causal tensors.

### 436 7.1. Ulam Map



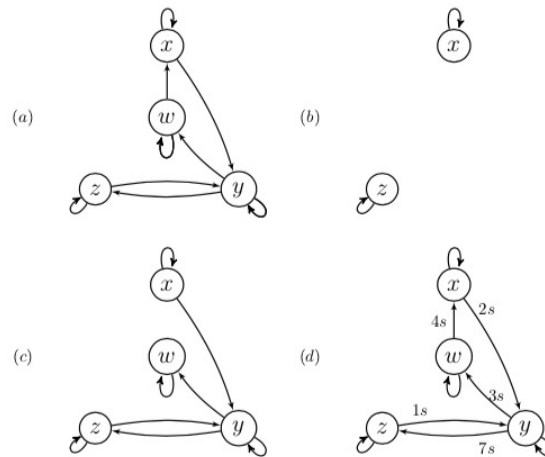
**Figure 4.** Transfer entropy and the channel capacity of the causal tensor for two unidirectionally coupled Ulam maps  $X^1$  and  $X^2$  as a function of the coupling strength  $\epsilon$ . Only the relation  $X^1 \rightarrow X^2$  is shown. Dots: approximated channel capacity for the causal channel. Line: transfer entropy as determined by Schreiber.

437 For the first experiment, we use the one-dimensional lattice of unidirectional coupled maps  
 438  $x_{n+1}^m = f(\epsilon x_n^{m-1} + (1 - \epsilon)x_n^m)$ . Information can only be transferred from  $X^{m-1}$  to  $X^m$ . The Ulam map  
 439 with  $f(x) = 2 - x^2$  is interesting because there are two regions ( $\epsilon \approx 0.18$ ,  $\epsilon \approx 0.82$ ) where no information  
 440 is shared between maps [8]. We chose an alphabet comprising four symbols. The quantization  
 441 comprised simple binning. Furthermore we chose  $\ell = m = 1$  (see Equation (3)). Instead of maximizing  
 442 TE we maximized the channel capacity to determine the optimal delay. An approximation that satisfies  
 443 the boundaries that follow from Equation (4) was used,

$$\tilde{\Gamma}(A) = \sum_g p(\psi_g^-) \Gamma(A_g). \quad (17)$$

444 To determine the channel capacities the Blahut-Arimoto algorithm was used [35]. The delays  
 445 were varied between one and 20. The Channel capacity was maximal for a delay of one sample. As  
 446 seen from Figure 4, causal tensors lead to a similar result as transfer entropy.

## 447 7.2. Coupled Ornstein-Uhlenbeck Processes



**Figure 5.** (a) The causal structure for the Ornstein-Uhlenbeck system of Equation (18). The other graphs show the inferred causal structures at different time series lengths. The confidence interval was 90% and the maximum delay was set to 20s: (b)  $T = 10ks$ , (c)  $T = 100ks$ . In (d),  $T = 500ks$ , the interaction delays that maximized the channel capacity are also shown.

In the second experiment, we demonstrate our approach using a system of four coupled Ornstein-Uhlenbeck processes [11]:

$$\begin{aligned}
 \dot{x}(t) &= -0.5x(t) + 0.6w(t-4)\eta_x(t), \\
 \dot{y}(t) &= -0.9y(t) - 1.0x(t-2) + 0.6z(t-5) + \eta_y(t), \\
 \dot{z}(t) &= -0.7z(t) - 0.5y(t-6) + \eta_z(t), \\
 \dot{w}(t) &= -0.8w(t) - 0.4y(t-3)^2 + 0.05y(t-3) + \eta_w(t),
 \end{aligned} \tag{18}$$

448 with independent unit variance white noise processes  $\eta$ . The integration time step was  $dt = 0.01s$   
 449 and the sampling interval  $\Delta s = 100s$ . We used a binary encoding scheme. First, the data was  
 450 normalized after which it was partitioned at 0.5. Because the Shannon entropy of the encoded data  
 451 was close to 1, we expect highly noisy communication channels. The disadvantage of binary encoding  
 452 is that more data is needed to capture the transmitted information. However, cascading very noisy  
 453 channels reduces the probability of detecting an indirect relation. This is illustrated by Figure 5, no  
 454 pruning was needed. This experiment shows that causal tensors can indeed detect the underlying  
 455 structure.

456 **8. Discussion**

457 In this article we focused mainly on the basic foundations of our formalism, we assumed for  
 458 example that there were no hidden causes. Whether our formalism can be used to detect the existence  
 459 of hidden causes will be researched in the future. If and how our framework applies to systems  
 460 comprising over three random variables will also be a subject of further research.

461 The conclusion that bivariate analysis suffices to differentiate direct and indirect associations,  
 462 follows directly from the DPI. It should not be of any surprise that in an information theoretical  
 463 framework, noise plays a fundamental role. A cascade of noisy channels is at least as noisy as the  
 464 noisiest channel within the cascade.



465 Further experiments are needed to confirm that a combination of multi-channel causal channels  
466 and single-channel causal channels could be used for pruning inferred networks.

467 The experiment in Subsection 7.1 illustrated that causal tensors give similar results as transfer  
468 entropy. The experiment in Subsection 7.2 illustrated that our formalism can indeed infer the  
469 underlying structure. All relations that should have been found based on the differential equations  
470 were found. All interaction delays except the interaction delay for  $Z \rightarrow Y$  were close or equal to  
471 the interaction delays in the differential equations. In both cases, we used a rather simple encoding  
472 scheme. This raises the question if and under what conditions an optimal encoding scheme exists  
473 that minimizes the cardinality of the alphabets, but keeps a sufficiently large amount of “causal  
474 information”.

475 We will focus future research on the extension and application of the framework. Because in the  
476 case of a directed triangle the tensor  $\bar{A}$  is related to the redundancy, the *potential* relation to “Partial  
477 Information Decomposition” [36] will be further explored in the future.

478 **Funding:** This research received no external funding.

479 **Acknowledgments:** I would like to thank Errol Zalmijn for introducing me to the wonderful topic of transfer  
480 entropy and Marcel Brunt for helping me to implement our approach in Matlab. Also, thanks to Hans Onvlee,  
481 S. Kolumban, Rui M. Castro and T. Heskes for their comments on earlier versions of the manuscript. ASML PI  
482 System Diagnostics supported part of the work.

483 **Conflicts of Interest:** The authors declare no conflict of interest.

## 484 Abbreviations

485 The following abbreviations are used in this manuscript:

486	DMC	Discrete memoryless communication channel
	MI	Mutual information
487	pmf	probability mass function
	TE	Transfer entropy

## 488 Appendix A. Proofs

### 489 Appendix A.1. Theorem 2

490 For the proof of Theorem 2 we need to introduce two lemma’s.

491 **Lemma A1.**  $\forall g : B_{gh\hat{j}}^k = B_{h\hat{j}}^k$ .

492 **Sketch of Proof for Lemma A1.** Another direct consequence of the Markov property is related to  
493 indices associated with the same random variable. As long as the index related to the past of the  
494 output— $g$ —and the index related to the output— $j$ —appear in the same tensor we are allowed to  
495 replace the output index by the input index. In our example, this means that we may replace  $j$  by  $\hat{j}$  as  
496 long as we ensure that  $\psi_{\hat{j}}^- = \{\psi_j, \psi_g^-\}$ . This is always possible because of the Markov property: we  
497 either enlarge the cardinality of  $\psi_{\hat{j}}^-$  or  $\psi_g^-$ .  $\square$

498 The next lemma follows directly from Example 1.

499 **Lemma A2.** For the chain  $X \rightarrow Y \rightarrow Z$  we have  $\mathcal{A}_{igh}^{\hat{i}} = \mathcal{A}_{ig}^{\hat{i}}$ .

500 **Sketch of Proof for Lemma A2.** Figure 1(a) depicts the situation of the chain  $X \rightarrow Y \rightarrow Z$ . According  
501 to the Causal Markov Condition  $\{\hat{i}, j\}$  and  $h$  are independent given  $g$ , i.e.,  $p(\hat{i}, j, h|g) = p(\hat{i}, j|g)p(h|g)$ .

502 We now rewrite this using expressions that follow from the definition for conditional probabilities.  
 503 The left-hand side is written as  $p(\hat{i}, j, h|g) = p(j, h|\hat{i}, g)p(\hat{i}|g)$ . The right-hand side can be rewritten  
 504 using  $p(\hat{i}, j|g) = p(j|\hat{i}, g)p(\hat{i}|g)$ .

505 This finally leads the conclusion that  $p(j, |\hat{i}, g, h) = p(j|\hat{i}, g)$ , i.e.,  $\mathcal{A}_{\hat{i}gh}^j = \mathcal{A}_{\hat{i}g}^j$ .  $\square$

506 **Sketch of Proof for Theorem 2.** Because of the Law of Total Probability we are allowed to condition  
 507 Equation (5) on  $h$  and both Equation (7a) and Equation (7b) on  $g$ . This leads to

$$p_{gh}^j = p_{gh}^{\hat{i}} A_{gh\hat{i}}^j \quad (\text{A1a})$$

$$p_{gh}^k = p_{gh}^{\hat{j}} B_{gh\hat{j}}^k \quad (\text{A1b})$$

$$p_{gh}^k = p_{gh}^{\hat{i}} C_{gh\hat{i}}^k \quad (\text{A1c})$$

508 Substituting the expression for  $p_{gh}^j$  of Equation (A1a) in Equation (A1b) and combining the result  
 509 with Equation (A1c) gives us  $C_{gh\hat{i}}^k = A_{gh\hat{i}}^j B_{gh\hat{j}}^k$ . Using Lemma(A1) and Lemma(A2) this can be rewritten  
 510 as

$$C_{gh\hat{i}}^k = A_{g\hat{i}}^j B_{h\hat{j}}^k \quad (\text{A2})$$

511 Finally, we multiply both sides with  $p_{h\hat{i}}^g$ . As the reader can confirm, the term  $p_{h\hat{i}}^g C_{gh\hat{i}}^k$  equals  $C_{h\hat{i}}^k$ .  
 512 This finally leads to Equation (8).

513

514 For the second part of the theorem, we refer to Figure 1(b). It depicts the situation of the directed  
 515 triangle  $X \rightarrow Y \rightarrow Z + X \rightarrow Z$ . According to the Causal Markov Condition  $\{\hat{i}, j\}$  and  $h$  are *not*  
 516 independent given  $g$ :  $p(\hat{i}, j, h|g) \neq p(\hat{i}, j|g)p(h|g)$ .

517 We now rewrite this using expressions that follow from the definition for conditional probabilities.  
 518 The left-hand side is written as  $p(\hat{i}, j, h|g) = p(j, h|\hat{i}, g)p(\hat{i}|g)$ . The right-hand side can be rewritten  
 519 using  $p(\hat{i}, j|g) = p(j|\hat{i}, g)p(\hat{i}|g)$ .  $\square$

#### 520 Appendix A.2. Lemma 1

521 **Sketch of Proof for Lemma 1.** By definition  $p_{h\hat{i}}^g A_{g\hat{i}}^j = \sum_g p(g|h, \hat{i})p(\hat{j}|g, \hat{i})$ . From Example 1(b) it  
 522 follows that  $\sum_g p(g|h, \hat{i})p(\hat{j}|g, \hat{i}) = \sum_g p(\hat{j}, g|\hat{i}, h)$ .

523 Applying the law of total probability to the righthand side gives us  $\sum_g p(\hat{j}, g|\hat{i}, h) = p(\hat{j}|\hat{i}, h)$ . In  
 524 other words:  $p_{h\hat{i}}^g A_{g\hat{i}}^j = p(\hat{j}|\hat{i}, h)$ .  $\square$

#### 525 Appendix A.3. Lemma 2

526 **Sketch of Proof for Lemma 2.** First we note that  $p_{h\hat{i}}^{\hat{j}\hat{j}} = \delta_{\hat{j}}^{\hat{j}} p_{h\hat{j}}^{\hat{j}} p_{h\hat{j}}^{\hat{i}}$ . Equation (4.3) can therefore be is  
 527 rewritten as  $p_h^k = \delta_{\hat{j}}^{\hat{j}} p_h^{\hat{j}} p_{h\hat{j}}^{\hat{i}} D_{h\hat{i}\hat{j}}^k$ . Because we are allowed to change the order of  $\delta_{\hat{j}}^{\hat{j}}$  and  $p_h^{\hat{j}}$  we get  $p_h^k =$

528  $p_h^{\hat{j}} \left( \delta_{\hat{j}}^{\hat{j}} p_{h\hat{j}}^{\hat{i}} D_{h\hat{i}\hat{j}}^k \right)$ . Combining this with Equation (7a) results in an expression for  $B_{h\hat{j}}^k$ :  $B_{h\hat{j}}^k = \delta_{\hat{j}}^{\hat{j}} p_{h\hat{j}}^{\hat{i}} D_{h\hat{i}\hat{j}}^k$ .

529 Because  $\delta_{\hat{j}}^{\hat{j}} p_{h\hat{j}}^{\hat{i}} = p_{h\hat{j}}^{\hat{i}}$  we get Equation (13a).  $\square$

#### 530 Appendix A.4. Theorem 6

531 For the proof of the data processing inequality, Theorem 6, the simplified notation for  
 532 transfer entropy and mutual information from Section 5 is used:  $TE_{X \rightarrow Y} := TE(\mathcal{A}, \cdot)$ ,  $TE_{Y \rightarrow Z} :=$

533  $TE(\mathcal{B}, \cdot)$ ,  $TE_{X \rightarrow Z} := TE(\mathcal{C}, \cdot)$  and  $I(X; Y) := I(\mathcal{A}_h, \cdot)$ ,  $I(Y; Z) := I(\mathcal{B}_h, \cdot)$ ,  $I(X; Z) := I(\mathcal{C}_h, \cdot)$ . The  
 534 subscript  $h$  indicates the  $h^{\text{th}}$  sub-channel representing a DMC.

535 **Sketch of Proof for Theorem 6.** We start with Equation (9) instead of Equation (14). The DPI is valid  
 536 per sub-channel. So, for all  $h$ :  $I(\mathcal{C}_h, \cdot) \leq \min[I(\mathcal{A}_h, \cdot), I(\mathcal{B}_h, \cdot)]$ . As per Equation (4) we multiply both  
 537 sides by  $p(\zeta_h^-)$ —the probability that the  $h^{\text{th}}$  channel is selected—and sum over  $h$ . This results in a DPI  
 538 for transfer entropy,

$$TE(\mathcal{C}, \cdot) \leq \min[TE(\bar{\mathcal{A}}, \cdot), TE(\mathcal{B}, \cdot)]. \quad (\text{A3})$$

The tensor  $\bar{\mathcal{A}}_h$  is itself the result of two cascaded channels represented by  $\mathcal{A}_g$  and a tensor with  
 elements  $p_{ih}^g$ . For these two DMC's the DPI is also valid, leading to:

$$\forall_{g,h} : I(\bar{\mathcal{A}}_h, \cdot) \leq I(\mathcal{A}_g, \cdot).$$

539 We now multiply both sides of this equation by  $p(\zeta_h^-)p(\psi_g^-)$ , and sum over  $h$  and  $g$ , resulting in  
 540  $TE(\bar{\mathcal{A}}, \cdot) \leq TE(\mathcal{A}, \cdot)$ . We can now rewrite Equation (A3) as

$$TE(\mathcal{C}, \cdot) \leq \min[TE(\mathcal{A}, \cdot), TE(\mathcal{B}, \cdot)]. \quad (\text{A4})$$

541 This implies that

$$TE(\mathcal{B} \odot \mathcal{A}, \cdot) \leq \min[TE(\mathcal{A}, \cdot), TE(\mathcal{B}, \cdot)]. \quad (\text{A5})$$

542  $\square$

#### 543 Appendix A.5. Theorem 7

544 **Sketch of the proof of Theorem 7.** If both  $\mathcal{B} = \mathcal{C} \odot \mathcal{A}^\ddagger$  and  $\mathcal{C} = \mathcal{B} \odot \mathcal{A}$  are valid, causal tensors can  
 545 not distinguish a fork from a chain. There are two cases that need to be considered. In the first case,  
 546 conditions are derived using the causal tensor relations. In the second case, we show that the pmfs  
 547 impose a certain condition.

548 We start by combining  $\mathcal{B} = \mathcal{C} \odot \mathcal{A}^\ddagger$  and  $\mathcal{C} = \mathcal{B} \odot \mathcal{A}$ :

$$\mathcal{B} = \mathcal{B} \odot \mathcal{A} \odot \mathcal{A}^\ddagger, \quad (\text{A6a})$$

$$\mathcal{C} = \mathcal{C} \odot \mathcal{A}^\ddagger \odot \mathcal{A}. \quad (\text{A6b})$$

549 These equations are valid when  $\mathcal{B} \odot \mathcal{I}_1 = \mathcal{B} \odot \mathcal{A} \odot \mathcal{A}^\ddagger$  and  $\mathcal{C} \odot \mathcal{I}_2 = \mathcal{C} \odot \mathcal{A}^\ddagger \odot \mathcal{A}$ , with  $\mathcal{I}_1$  and  $\mathcal{I}_2$   
 550 identity causal tensors. Per definition identity tensors are noiseless. Because the causal tensors are  
 551 stochastic tensors, their elements are nonnegative. The product of two stochastic tensors can only  
 552 equal a noiseless tensor if and only if both  $\bar{\mathcal{A}}$  and  $\bar{\mathcal{A}}^\ddagger$  are noiseless. Along the same line of reasoning,  
 553 we finally conclude that  $\mathcal{A}$  and  $\mathcal{A}^\ddagger$  are noiseless causal tensors because the averaging operation is in  
 554 fact a matrix multiplication of two tensors.

555

556 The second case in which we cannot distinguish a fork and a chain follows from the pmf  
 557 transformations:

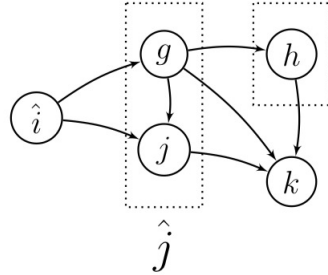
$$\mathcal{B} | y \rangle = \mathcal{C} \odot \mathcal{A}^\ddagger | y \rangle, \quad (\text{A7a})$$

$$\mathcal{C} | x \rangle = \mathcal{B} \odot \mathcal{A} | x \rangle. \quad (\text{A7b})$$

558 The output from both the left-hand side and the right-hand side of these equations are probability  
 559 mass functions. If they are indistinguishable, we cannot differentiate between a fork and a chain either.

560 If both  $\mathcal{B}$  and  $\mathcal{C}$  are (almost) perfectly noisy causal tensors, both equations in Equation (A7) reduce to  
 561  $|u\rangle = |u\rangle$ , with  $|u\rangle$  representing the uniform pmf. We cannot distinguish a chain from a fork.  $\square$

562 *Appendix A.6. Theorem 8*



**Figure A1.** The graph related a chain of two causal channels. Both dotted boxes block the paths between  $\hat{i}$  and  $k$ . The dotted box comprising the variables  $\{g, j\}$  also blocks all paths between  $\hat{i}$  and  $k$ . This box represents the variable  $\hat{j}$ .

563 The proof for Theorem 8 uses Figure A1. We furthermore do not use the Einstein summation  
 564 convention.

565 **Sketch of the proof of Theorem 8.** According to the Law of Total Probability,  $C_i^k = \sum_h p_i^h C_{hi}^k$ .  
 566 Multiplying both sides of Equation (9) by  $p_i^h$ , gives us

$$C_i^k = \sum_j \sum_h p_i^h A_{hi}^j B_{hj}^k. \quad (\text{A8})$$

We now express all stochastic tensors with the letter  $p$  instead of an  $A$  or a  $B$ . Using the fact that  $p_i^h p_{hi}^j = p_i^j$ , and  $p_{hj}^k = p_j^{kh} / p_j^h$ , Equation (A8) can be written as:

$$C_i^k = \sum_j \sum_h p_i^j p_j^{kh} / p_j^h. \quad (\text{A9})$$

567 Because  $p_i^j = p_i^h p_{ij}^h$  and  $p_j^{kh} = p_j^k p_{jk}^h$  we get

$$C_i^k = \sum_j p_i^j p_j^k \sum_h p_{ij}^h p_{jk}^h / p_j^h. \quad (\text{A10})$$

568 Using  $p_{ij}^h = p_j^{hi} / p_j^i$  and  $p_{jk}^h = p_j^{hk} / p_j^k$  we rewrite Equation (A10) as

$$C_i^k = \sum_j p_i^j p_j^k \sum_h \frac{p_j^{hi}}{p_j^i} \frac{p_j^{hk}}{p_j^k} \frac{1}{p_j^h}. \quad (\text{A11})$$

Using a similar step as in the previous cases, Equation (A11) can be rewritten as

$$C_i^k = \sum_j p_i^j p_j^k \sum_h \frac{p_{hj}^i}{p_j^i} \frac{p_{hj}^k}{p_j^k} p_j^h. \quad (\text{A12})$$

569 Applying the Causal Markov Condition on the graph in Figure A1, gives us  $\hat{i} \perp\!\!\!\perp k | \{h, \hat{j}\}$ , i.e.,  
 570  $p_{hj}^i p_{hj}^k p_j^h = p_j^{hk}$ . Because we sum over  $h$ , we sum out this variable as per Law of Total Probability. Using  
 571  $A_i^j = p_j^j$  and  $B_j^k = p_j^k$  we finally get the following expression

$$C_i^k = \sum_j A_i^j B_j^k \frac{p_j^{\hat{i}k}}{p_j^{\hat{i}} p_j^k}. \quad (\text{A13})$$

572 So, if  $\hat{i}$  and  $k$  are independent given  $\hat{j}$ , the theorem has been proven. When we apply the Causal  
573 Markov Condition to the graph in Figure A1, we see that  $\hat{i} \perp\!\!\!\perp k | \hat{j}$ .  $\square$

## 574 References

- 575 1. Guo, R.; Cheng, L.; Li, J.; Hahn, P.; Liu, H. A Survey of Learning Causality with Data: Problems and  
576 Methods **2018**.
- 577 2. Eichler, M. Causal inference with multiple time series: Principles and problems. *Philosophical transactions.*  
578 *Series A, Mathematical, physical, and engineering sciences* **2013**, 371, 20110613. doi:10.1098/rsta.2011.0613.
- 579 3. Pearl, J. *Causality: Models, Reasoning and Inference*, 2nd ed.; Cambridge University Press: New York, NY,  
580 USA, 2009.
- 581 4. Granger, C.W.J. Investigating Causal Relations by Econometric Models and Cross-Spectral Methods.  
582 *Econometrica* **1969**, 37, 424–438.
- 583 5. Vastano, J.A.; Swinney, H.L. Information transport in spatiotemporal systems. *Phys. Rev. Lett.* **1988**,  
584 60, 1773–1776. doi:10.1103/PhysRevLett.60.1773.
- 585 6. Dagum, P.; Galper, A.; Horvitz, E. Dynamic Network Models for Forecasting. Proceedings of the Eighth  
586 International Conference on Uncertainty in Artificial Intelligence; Morgan Kaufmann Publishers Inc.: San  
587 Francisco, CA, USA, 1992; UAI'92, pp. 41–48.
- 588 7. Spirtes, P.; Glymour, C.; N., S.; Richard. *Causation, Prediction, and Search*; Mit Press: Cambridge, 2000.
- 589 8. Schreiber, T. Measuring Information Transfer. *Phys. Rev. Lett.* **2000**, 85, 461–464.  
590 doi:10.1103/PhysRevLett.85.461.
- 591 9. Lizier, J.T.; Prokopenko, M. Differentiating information transfer and causal effect. *The European Physical*  
592 *Journal B* **2010**, 73, 605–615. doi:10.1140/epjb/e2010-00034-5.
- 593 10. Hyvärinen, A.; Zhang, K.; Shimizu, S.; Hoyer, P.O. Estimation of a Structural Vector Autoregression Model  
594 Using Non-Gaussianity. *J. Mach. Learn. Res.* **2010**, 11, 1709–1731.
- 595 11. Runge, J.; Heitzig, J.; Petoukhov, V.; Kurths, J. Escaping the Curse of Dimensionality in Estimating  
596 Multivariate Transfer Entropy. *Phys. Rev. Lett.* **2012**, 108, 258701. doi:10.1103/PhysRevLett.108.258701.
- 597 12. Duan, P.; Yang, F.; Chen, T.; Shah, S. Direct Causality Detection via the Transfer Entropy Approach. *Control*  
598 *Systems Technology, IEEE Transactions on* **2013**, 21, 2052–2066. doi:10.1109/TCST.2012.2233476.
- 599 13. Sun, J.; Taylor, D.; Bollt, E. Causal Network Inference by Optimal Causation Entropy. *SIAM Journal on*  
600 *Applied Dynamical Systems* **2014**, 14. doi:10.1137/140956166.
- 601 14. Turing, A.M. On Computable Numbers, with an Application to the  
602 Entscheidungsproblem. *Proceedings of the London Mathematical Society* **1937**,  
603 s2-42, 230–265, [<http://oup.prod.sis.lan/plms/article-pdf/s2-42/1/230/4317544/s2-42-1-230.pdf>].  
604 doi:10.1112/plms/s2-42.1.230.
- 605 15. Copeland, B.J. The Church-Turing Thesis. In *The Stanford Encyclopedia of Philosophy*, Spring 2019 ed.; Zalta,  
606 E.N., Ed.; Metaphysics Research Lab, Stanford University, 2019.
- 607 16. Pearl, J. Causal Inference in Statistics: An Overview. *Statistics Surveys* **2009**, 3, 96–146. doi:10.1214/09-SS057.
- 608 17. Eberhardt, F.; Scheines, R. Interventions and Causal Inference. *Philos. Sci.* **2007**, 74. doi:10.1086/525638.
- 609 18. Ding, M.; Chen, Y.; Bressler, S.L. Granger Causality: Basic Theory and Application to Neuroscience, 2006.
- 610 19. Ghassami, A.; Kiyavash, N. Interaction information for causal inference: The case of directed  
611 triangle. 2017 IEEE International Symposium on Information Theory (ISIT), 2017, pp. 1326–1330.  
612 doi:10.1109/ISIT.2017.8006744.
- 613 20. Cover, T.M.; Thomas, J.A. *Elements of Information Theory*; Wiley-Interscience: New York, NY, USA, 1991.
- 614 21. Margolin, A.; Nemenman, I.; Basso, K.; Wiggins, C.; Stolovitzky, G.; Dalla-Favera, R.; Califano, A. ARACNE:  
615 An Algorithm for the Reconstruction of Gene Regulatory Networks in a Mammalian Cellular Context.  
616 *BMC bioinformatics* **2006**, 7 Suppl 1, S7. doi:10.1186/1471-2105-7-S1-S7.
- 617 22. Papoulis, A.; Pillai, S.U. *Probability, Random Variables, and Stochastic Processes*, fourth ed.; McGraw Hill,  
618 2002.

- 619 23. Shannon, C.E. A Mathematical Theory of Communication. *Bell System Technical Journal*,  
620 27, 379–423, [<https://onlinelibrary.wiley.com/doi/pdf/10.1002/j.1538-7305.1948.tb01338.x>].  
621 doi:10.1002/j.1538-7305.1948.tb01338.x.
- 622 24. Ahmed, S.S.; Roy, S.; Kalita, J.K. Assessing the Effectiveness of Causality Inference Methods for Gene  
623 Regulatory Networks. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **2018**, pp. 1–1.  
624 doi:10.1109/TCBB.2018.2853728.
- 625 25. Rashidi, B.; Singh, D.S.; Zhao, Q. Data-driven root-cause fault diagnosis for  
626 multivariate non-linear processes. *Control Engineering Practice* **2018**, *70*, 134 – 147.  
627 doi:<https://doi.org/10.1016/j.conengprac.2017.09.021>.
- 628 26. Dullemond, K.; Peeters, K. Introduction to Tensor Calculus. 2010.
- 629 27. Muroga, S. On the Capacity of a Discrete Channel. *Journal of the Physical Society of Japan* **1953**, *8*, 484–494.  
630 doi:10.1143/JPSJ.8.484.
- 631 28. Beckenbach, E.F. *Modern mathematics for the engineer: second series*; New York : McGraw-Hill, 1961.
- 632 29. Wibral, M.; Pampu, N.; Priesemann, V.; Siebenhühner, F.; Seiwert, H.; Lindner, M.; Lizier, J.T.; Vicente, R.  
633 Measuring information-transfer delays. *PloS one* **2013**.
- 634 30. Dean, T. *Network+ Guide to Networks*, 6th ed.; Course Technology Press: Boston, MA, United States, 2012.
- 635 31. Berge, C. *Graphs and Hypergraphs*; Elsevier Science Ltd.: Oxford, UK, UK, 1985.
- 636 32. Dirac, P.A.M. A new notation for quantum mechanics. *Mathematical Proceedings of the Cambridge Philosophical*  
637 *Society* **1939**, *35*, 416–418. doi:10.1017/S0305004100021162.
- 638 33. Nauta, M.; Bucur, D.; Seifert, C. Causal Discovery with Attention-Based Convolutional Neural Networks.  
639 *Machine Learning and Knowledge Extraction* **2019**, *1*, 312–340. doi:10.3390/make1010019.
- 640 34. Wibral, M.; Wollstadt, P.; Meyer, U.; Pampu, N.; Priesemann, V.; Vicente, R. Revisiting Wiener’s principle  
641 of causality — interaction-delay reconstruction using transfer entropy and multivariate analysis on  
642 delay-weighted graphs. 2012 Annual International Conference of the IEEE Engineering in Medicine  
643 and Biology Society, 2012, pp. 3676–3679. doi:10.1109/EMBC.2012.6346764.
- 644 35. Blahut, R. Computation of channel capacity and rate-distortion functions. *IEEE Transactions on Information*  
645 *Theory* **1972**, *18*, 460–473. doi:10.1109/TIT.1972.1054855.
- 646 36. Williams, P.; Beer, R. Nonnegative Decomposition of Multivariate Information. *preprint* **2010**, 1004.