

Information Flow in Computational Systems

Praveen Venkatesh*, Sanghamitra Dutta† and Pulkit Grover‡

Electrical & Computer Engineering, and the Center for the Neural Basis of Cognition
Carnegie Mellon University

*vpraveen@cmu.edu †sanghamd@andrew.cmu.edu ‡pulkit@cmu.edu

March 3, 2020

Abstract

We develop a theoretical framework for defining and identifying flows of information in computational systems. Here, a computational system is assumed to be a directed graph, with “clocked” nodes that send transmissions to each other along the edges of the graph at discrete points in time. We are interested in a definition that captures the dynamic flow of information about a specific *message*, and which guarantees an unbroken “information path” between appropriately defined inputs and outputs in the directed graph. Prior measures, including those based on Granger Causality and Directed Information, fail to provide clear assumptions and guarantees about when they correctly reflect information flow about a message. We take a systematic approach—iterating through candidate definitions and counterexamples—to arrive at a definition for information flow that is based on conditional mutual information, and which satisfies desirable properties, including the existence of information paths. Finally, we describe how information flow might be detected in a noiseless setting, and provide an algorithm to identify information paths on the time-unrolled graph of a computational system.

1 Introduction

1.1 Motivation

Neuroscientists¹ often seek an understanding of how information flows in the brain while it performs a particular task [2–5]. As a concrete example, consider the experiment performed by Almeida et al. [2], where they examine how images of common handheld tools are processed in the brain. In simple terms, the question they investigate is this: when attempting to identify a handheld tool, does one make use of knowledge of how to manipulate it? Two hypotheses present themselves: (i) the answer to the above question is *yes*, so we should expect that information about a tool’s identity *first* flows from visual cortex to motor cortex (the area responsible for processing manipulation), *before* synthesis of visual and motor information occurs at the area of the brain responsible for object recognition. (ii) The answer to the aforementioned question is *no*, so we should expect that the information about tools’ identities *first* flows from visual cortex to the area responsible for object recognition, *after* which this information arrives at motor cortex. Thus, distinguishing between these hypotheses is equivalent to determining the *path* along which information about a tool’s identity flows in the brain. What methods can neuroscientists use to gain such an understanding? What formal theory underlies such an analysis? How does one mathematically define colloquially-used terms such as “information flow”? These are the fundamental questions we try to answer in this paper.

¹A short version of this paper has appeared in the 2019 IEEE International Symposium on Information Theory [1].

Information flow is a concept that appears in several contexts, across fields ranging from communication systems, control theory and neuroscience to security, algorithmic transparency, and deep learning. While our primary motivation comes from neuroscience, the theory that we develop is broadly applicable to any system which can be modeled in the form of a directed graph, with nodes that communicate functions of their inputs to other nodes, and where transmissions are observable. For example, several kinds of social networks readily fit this bill, and one might wish to analyze how information spreads in such networks. Our framework is also general enough to analyze information flow in various kinds of Artificial Neural Networks: this could be useful for identifying specific paths that carry information distinguishing two or more classes, or for intelligently pruning an Artificial Neural Network post-training.

In the field of neuroscience, studying normal and diseased brain function involves gaining insight into how information is processed in the brain. Attaining such insight, in turn, may require determining how information flows between various parts of the brain. Thus, a nuanced understanding of information flow in the brain could help with diagnosing and treating brain diseases: a subject that is currently of immense interest with numerous efforts around the world [6–9]. More generally, an understanding of information flow is essential when considering how one might intervene to affect the output of a computational system, be it modulating how information spreads in a social network, or complementing dysfunctional components of the nervous system through stimulation (such as in retinal and cochlear implants).

1.2 Our Goal and Approach

Our overarching goal in this paper is develop a formal theory for understanding information flow in neuroscientific experiments. In order to properly scope our task, we choose to restrict our attention to “*event related experimental paradigms*” [10], a set of standard neuroscientific experimental design protocols where a *stimulus* is shown to an animal subject or human participant, whose brain signals are being recorded. This restriction also allows us to decide precisely what *kind* of information flow we are interested in, since in general, the phrase “information flow” can refer to more than one notion in neuroscience. We identify two dominant interpretations of “information flow”: (i) the first refers to information about a specific quantity or variable that is of interest to the experimentalist, which in this paper we refer to as the “message”; (ii) the second refers to information in the abstract, and is usually used to describe the fact that one area of the brain “drives” or “influences” another area through the transmission of some information: in this interpretation, one is not interested in *what* is being communicated, only that the communication is *occurring*. In this paper, we focus only on the first interpretation of the phrase, where we are interested in information about a *specific message*, and we wish to track how information about this message flows within the brain. All references to “information flow”, henceforth, refer only to the first interpretation. This is particularly common in event-driven paradigms, where the neuroscientist investigates how the brain responds to a carefully chosen set of stimuli, and examines how information contained in these stimuli (or alternatively, information contained in the *response*) flows through the brain.

Given that we are interested in information about a specific message, what are we in pursuit of when we say information *flow*? Broadly speaking, we want to develop a measure that will allow us to examine *how* and *at what times* information about a specific message flows from one area of the brain to another. In particular, we think of the brain as a *computational system* executing an algorithm, and we want to capture how information about different variables might flow between different computational nodes of this system. A given “message” variable may be stored at a particular node for some time, a function may be computed using this message, and then the result may be passed on to a different node. A node that transmits information about the message at one time instant may not do so at a later time instant; thus information flow is a dynamic or time-dependent quantity. The computational system should allow for all of these possibilities,

and our ensuing measure of information flow should enable us to *track the path traversed by the message² through this system, over time*. This is the principal goal of our theoretical development, and will guide many of our decisions in model design.

We approach this goal by formally defining a computational system model: one based on nodes that represent distinct computational areas of the brain. These nodes can potentially represent the brain at any scale: single neurons, groups of neurons, or even whole brain regions, depending on the measurement modality and the kind of experiment being performed. The computational model we develop borrows various ideas from across several fields. The basic ingredients of the computational model, based on a *graph with computational nodes*, derives from Thompson’s work on VLSI complexity theory [11]. In order to attain a *dynamic* picture of information flow on the edges of the computational graph, and to deal with cycles in the flow of information, we use the idea of time-unrolling a graph, taking inspiration from Network Information Theory [12]. Finally, to describe how the computational nodes can compute stochastic functions of variables based on their current inputs, we use the idea of Structural Causal Models from the field of Causality [13, 14].

Within this computational model, we define a new measure for information flow about a *specific message* that captures the *dynamic* nature of information transmission. Ultimately, the measure should allow us to track how information about the message flows through the system, in the form of an *unbroken information path*—this will be a key focus guiding our definitions. Given the nature of the problem, we rely on information-theoretic measures to define information flow. We motivate this definition through properties, and provide a series of candidate definitions and counterexamples before arriving at our final definition. When defining information flow in such a computational system, we restrict ourselves to “*observational*” measures, which can be computed from a sample of all random variables described in the model. We deliberately eschew *interventional* and *counterfactual* measures, as the former require the capability to intervene on the system and change the distributions of the random variables involved, while the latter are a purely theoretical notion that can only be applied in situations where one can ask what *might have occurred* if a specific variable had been different on a particular trial (while keeping the realizations of all other latent sources of randomness fixed).

The approach of building a rigorous theoretical framework that we have adopted in this paper is inspired by two works from biologists titled “Can a biologist fix a radio?” [15] and “Could a neuroscientist understand a microprocessor?” [16]. Both these works point to the lack of formal methods, i.e., systematic theory, that could help biologists understand the limitations of their tools and test their assumptions. It is our belief that information theory can help provide the formal methods that are sought in biology, and make an impact in fields such as neuroscience and neuroengineering [17–19]. In particular, information theory can play an important role in advancing how we understand large computational systems through external measurements and interventions. While developing an understanding of information flow in such systems may not be sufficient for providing a complete description of the nature of computation itself, we believe that it forms an integral component. Going forward, we believe that providing a formal theoretical framework for information flow is but a small part of several larger theoretical questions that are yet to be properly posed: questions such as how one might formalize “reverse engineering” the brain, or formalize the notion of “understanding computation”.

1.3 Related Work

Prior work on statistically inferring flows of information in the brain appears under the umbrella of “functional” or “effective connectivity” [20–22]. These efforts have largely relied on measures of statistical³ causal

²or information derived thereof

³We borrow the use of the term “statistical” from Pearl [13, Sec. 1.5], who contrasts and differentiates “statistical” concepts from (strictly) “causal” ones.

influence such as Granger Causality [23, 24], Massey’s Directed Information [25–28], Transfer Entropy [29] and Partial Directed Coherence [30]. Despite widespread use, these measures have frequently been a subject of debate and disagreement within the neuroscientific community [31–37]. In part, these disagreements stem from the widely-acknowledged fact that under non-ideal measurement conditions (e.g. in the presence of hidden variables [13, p. 54], asymmetric noise [38, 39], or limited sampling [40]), estimation of these quantities may be erroneous. While these non-idealities may eventually be overcome through improvements in technology, we believe that more fundamental issues still remain. For instance, one basic question that has remained unanswered is: when can statistical causal influence be interpreted as information flow about a message? In previous work, we demonstrated that even under *ideal* measurement conditions, the direction of greater Granger causal influence can be opposite to the direction in which the message is being communicated in certain kinds of feedback communication networks [41]. This example points to a more general issue with the use of statistical causal influence measures: there is no direct way to interpret what the influence is “*about*”. While it is understood in certain settings that “information flow” refers to information contained in a particular set of “*stimuli*” (as mentioned in the previous section), the aforementioned measures do not incorporate the effect of the stimulus.

The existence of such fundamental issues can be traced back to the fact that there is no underlying model that links information flow (of some message of interest) with the signals that are actually *measured*, leading to a lack of separation between the problems of *defining* information flow and of *estimating* it. The lack of such a computational model also makes it hard to test assumptions and to draw the right interpretations from experimental analyses. We believe that, following Shannon’s approach of providing a theoretical foundation for information transmission [42], a solid theoretical treatment of information flow is needed. Such a treatment would begin with a model of the underlying system, give a definition for information flow and describe its properties, and finally end with a suitable estimator. Adopting Shannon’s model of defining entropy by stating a set of properties that such a measure must satisfy, we attempt to define information flow by putting forward an intuitive property that we believe is desirable for such a quantity. It is our hope that, by providing a theoretical foundation that separates definition and estimation, along with a concrete model and explicitly-stated assumptions, we can avoid many of the pitfalls encountered by previous approaches to understanding information flow in the brain.

It is useful at this point to mention the key differences between our measure of information flow, and measures based on Granger Causality and its generalizations:

1. Our measure depends on a message M , that will often be related to the stimulus or the response in a neuroscientific task, whereas tools based on Granger causality do not.
2. Since Granger causality-based tools use time series modeling to compute an estimate of information flow, they are unable to provide a dynamic, evolving picture of information flow between different areas over time.
3. Since we start with a *computational framework*, our model provides a direct way to connect information flow with the underlying computation. On the other hand, Granger causality-based tools start with a probabilistic graphical model of the observed nodes, and do not tie the analysis to computation in any way.

While our proposed definition of information flow will also suffer from performance degradation under non-ideal measurement conditions, we believe that it overcomes the fundamental difficulty faced by Granger Causality-based tools: when measurements are ideal, our definition provides a clear and consistent way to interpret information flow about a message, as we illustrate through several examples in Section 6.

Another line of work that appears within the functional and effective connectivity literature is Dynamic Causal Modeling (DCM) [21, 43]. This methodology is, in spirit, much more closely aligned with what

we propose here. However, our framework differs from DCM in a few important ways: (i) our underlying framework and model is based on Structural Causal Models rather than dynamical systems, and (ii) we seek to formalize the notion of information flow, not just of effective connectivity. However, the style of thinking, which involves starting from theoretical models and incorporating the stimulus and experimental design, is common to both DCM and our approach.

1.4 Outline of the Paper

In this paper, we start by giving a mathematical description of a generic computational system, about which inferences are being drawn (Section 2). We then formally define what it means for information about some message to flow on a single edge or on a set of edges in the computational system (Section 3). This is done by proposing an intuitive property that we would like such flows to satisfy, along with some candidate definitions, and then examining which candidates satisfy the property. The intuitive property we desire is: *information flow about a message may not completely disappear from the system at a certain time, only to spontaneously reappear at a later point* (formalized in Property 1). It emerges that simple and intuitive definitions actually fail to satisfy this basic property, and so a more sophisticated definition is needed. We then show how our definition for information flow about the message satisfies several desirable properties, including guarantees for the existence of so-called “information paths” between appropriately defined input and output nodes (Section 4). After this, we suggest how one might detect which edges of the computational system have information flow, and provide an “information path algorithm”, which identifies the aforementioned information paths (Section 5). We also introduce and discuss the concepts of derived information, redundant transmissions and hidden nodes, which allow one to obtain a more fine-grained understanding of information structure in the computational system. To show that our definition of information flow agrees with intuition, we give several canonical examples of computational systems and depict the information flow in each case (Section 6). Finally, we conclude with discussions on connections with neuroscience, issues related to the difficulty of estimating information flow (along with possible remedies), comparison with the existing directed causal influence literature, connections with fields such as probabilistic graphical models and causality, and a discussion on information volume (Section 7).

2 The Computational System

Our goal is to develop a rigorous framework for understanding how the information about a message flows in a computational system. To do this, we first need to define the terms “computational system”, “message”, “information about a message” and “flow”. In this section, we start with the first two terms, defining the model of the computational system that is used throughout this paper, and explicitly defining the message.

Our model is based on prior art in the information theory literature [11, 12], and consists of nodes communicating to each other at discrete points in time on a directed graph. At every time instant, each node receives transmissions on its incoming edges and computes a function of these transmissions to send out on its outgoing edges. This function can be random and time-dependent, and can be different for every outgoing edge. We will be interested in the flow of a particular random variable called the “message”, which will be defined shortly. Since the directed graph forming the computational system may have cycles, the message may flow along a cyclic path. To deal with this possibility while capturing the fact that nodes must be causal⁴, we define a “time-unrolled” graph (in a manner similar to Ahlswede et al.⁵ [12]), which

⁴Causal in the “Signals and Systems” sense of the word, where a node cannot make use of future transmissions [44].

⁵Although the work of Ahlswede et al. (2000) is titled “Network *Information Flow*”, it actually addresses a different problem: one of the achievable rate region of a broadcast network and the optimal coding strategy that achieves this rate. In contrast to their work, which concentrates on characterizing and achieving the optimal rate, our focus is on understanding how information

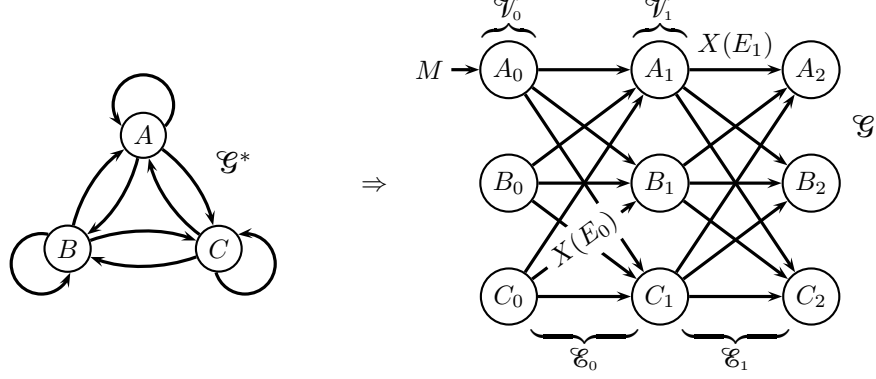


Figure 1: A diagram showing an example of how a complete directed graph is unrolled to create a time-unrolled graph. On the left, we show a complete directed graph \mathcal{G}^* that has three nodes, $\mathcal{V}^* = \{A, B, C\}$. These nodes are fully connected to each other via edges \mathcal{E}^* , including self-edges.

On the right, we show how \mathcal{G}^* has been unrolled using time indices $\mathcal{T} = \{0, 1, 2\}$ to obtain a time-unrolled graph \mathcal{G} . The set of all nodes at time $t = 0$ is \mathcal{V}_0 and the set of all (outgoing) edges at time $t = 0$ is denoted \mathcal{E}_0 . As an example, we have shown an arbitrary edge $E_0 \in \mathcal{E}_0$ (here, $E_0 = (C_0, B_1)$) and the transmission on that edge, $X(E_0)$. As another example, we show a “self-edge” in the time-unrolled graph, $E_1 \in \mathcal{E}_1$, which in this case is $E_1 = (A_1, A_2)$. Also depicted is the transmission $X(E_1)$ on this self-edge, which is interpreted as the contents of the memory of node A from $t = 1$ to $t = 2$. The message M arrives at the input node A_0 , but could in general be available at more than one node at $t = 0$.

In subsequent illustrations, we do not depict all edges at every time step, even though they are present. This is done only for the sake of clarity.

describes how nodes communicate to each other over time. We define a random variable model for the nodes’ transmissions, and demonstrate how each node computes these variables. We also formally define the input nodes of the computational system, through their relationship with the message.

Definition 1 (Complete directed graph). *A complete directed graph $\mathcal{G}^* = (\mathcal{V}^*, \mathcal{E}^*)$ is described by a set of nodes and the set of all edges between those nodes (including self-edges). We denote the set of nodes by their indices, $\mathcal{V}^* = \{1, 2, \dots, N\}$, where N is a positive integer denoting the number of nodes in the graph. The set of edges in the graph is the set of all ordered pairs of nodes, $\mathcal{E}^* = \mathcal{V}^* \times \mathcal{V}^*$.*

Note that (i) edges are directed, so the edge $(A, B) \in \mathcal{E}^*$ describes an edge *from* node A to node B ; and (ii) nodes have self-edges. For every $A \in \mathcal{V}^*$, there is an edge (A, A) in \mathcal{E}^* .

Moving forward, nodes shall be thought of as performing computations and possessing local memories. We shall interpret the transmission of a node to itself as the variable it stores within its memory⁶.

Definition 2 (Time-unrolled graph). *In order to allow nodes to have different transmissions at every time instant, we must provide for the progression of time. Let $\mathcal{T} = \{0, 1, \dots, T\}$ be a set of time indices, where T is a positive integer representing the maximum time index. Then, a time-unrolled graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is constructed by indexing a complete directed graph \mathcal{G}^* using the time indices \mathcal{T} as follows:*

1. The nodes \mathcal{V} consist of all nodes \mathcal{V}^* in \mathcal{G}^* , subscripted by time indices in \mathcal{T} ,

$$\mathcal{V} = \{A_t : A \in \mathcal{V}^*, t \in \mathcal{T}\};$$

2. The edges \mathcal{E} connect nodes of successive times in \mathcal{V} , so they can be written in terms of the edges in \mathcal{E}^* as

$$\mathcal{E} = \{(A_t, B_{t+1}) : (A, B) \in \mathcal{E}^*, t \in \mathcal{T}\}.$$

about a known message flows in an existing computational system.

⁶Instances of directed graphs that are *not complete* and of nodes possessing *no memory* are merely special cases of our model, where the respective edges’ transmissions can simply be set to zero.

For brevity, we denote the set of all nodes at time t by \mathcal{V}_t , and the set of all (outgoing) edges at time t by \mathcal{E}_t . So, for example, we will have $A_1 \in \mathcal{V}_1$ and $(A_1, B_2) \in \mathcal{E}_1$. All of the notation in this section can be visualized in Figure 1 and is summarized in Table 1.

Once again, note that (i) edges at time t connect nodes at time t to nodes at time $t + 1$; and (ii) since the original graph \mathcal{G}^* had self-edges, there will always be an edge (A_t, A_{t+1}) in \mathcal{E}_t for every node $A_t \in \mathcal{V}_t$.

Also note, we have only presented the complete directed graph in Definition 1 in order to explicitly define the process of time-unrolling. We do not expect the time-unrolled graph to be “rolled back” into a complete directed graph at the end of an information flow analysis. Since we seek a time-evolving picture of information flow between different computational nodes, we will directly view and interpret information flow on the time-unrolled graph. This is illustrated later, through several examples, in Section 6.

Definition 3 (Computational System). *A computational system $\mathcal{C} = (\mathcal{G}, X, W, f)$ is a time unrolled graph \mathcal{G} that has transmissions on its edges which are constrained by computations at its nodes. The input to the computational system includes a message⁷, M . We now elaborate upon these terms:*

3a) Transmissions on Edges

We begin by defining a function which maps every edge of \mathcal{G} to a random variable. Let \mathcal{X} be the set of all random variables in some probability space⁸. Then, let $X : \mathcal{E} \rightarrow \mathcal{X}$ be a function that describes what random variable is being transmitted on a given edge, i.e., $X(E)$ is the random variable corresponding to the transmission on the edge E .

For convenience, we define X applied to a set of edges as the set of random variables produced by applying X to each of those edges individually, i.e., for any set $\mathcal{E}' \subseteq \mathcal{E}$,

$$X(\mathcal{E}') = \{X(E) : E \in \mathcal{E}'\}. \quad (1)$$

We extend the use of this notation to other functions of nodes and edges that we define, going forward.

3b) Computation at a Node

Let $A_t \in \mathcal{V}_t$ be a node in the time-unrolled graph \mathcal{G} , at some time $t \geq 1$ (recall that $t \in \{0, 1, \dots, T\}$). Let $\mathcal{P}(A_t)$ be the set of edges entering A_t , and $\mathcal{Q}(A_t)$ be the set of edges leaving A_t . Further, let us suppose that A_t is able to intrinsically generate the random variable⁹ $W(A_t)$ at time t , where $W(A_t) \perp W(\mathcal{V} \setminus \{A_t\}) \forall A_t \in \mathcal{V}$, $W(\mathcal{V}_t) \perp \{M, X(\mathcal{E}_{t-1})\}$ and the symbol “ \perp ” stands for independence between random variables. Then, the computation performed by the node A_t (for $t \geq 1$) is a deterministic function¹⁰ f_{A_t} that satisfies

$$f_{A_t}(X(\mathcal{P}(A_t)), W(A_t)) = X(\mathcal{Q}(A_t)). \quad (2)$$

Here, $X(\mathcal{E}_{t-1})$, $W(\mathcal{V} \setminus \{A_t\})$, $W(\mathcal{V}_t)$, $X(\mathcal{P}(A_t))$ and $X(\mathcal{Q}(A_t))$ all make use of the notation described in (1).

Note that the definition above does not apply when $t = 0$; this is a special case which is discussed below. Also, for convenience, where \mathcal{A} is an arbitrary set of nodes, we will use $f_{\mathcal{A}}$ to denote the “joint function” mapping the incoming transmissions of all nodes in \mathcal{A} (along with their intrinsic random variables $W(\mathcal{A})$) to their respective outgoing transmissions.

⁷The message is the random variable whose “information flow” we will seek to identify.

⁸We assume that all probability distributions are such that the mutual information and conditional mutual information between any sets of random variables is well-defined [45, Sec. 2.6].

⁹ $X(E_t)$ and $W(A_t)$ may also be random *vectors* instead of random variables, i.e., an edge may *transmit a vector*. This does not affect the theoretical development presented in this paper; all of our proofs remain unchanged.

¹⁰This kind of model is not new, and can be found in the causality literature for instance, under the name “Structural Equation Models” [13, Sec. 1.4.1].

3c) The Message and the Input Nodes

Each of the nodes in \mathcal{V}_0 may receive one or more random variables from the world external to the computational system at time $t = 0$. The message, M , is simply a specific random variable that is of interest to the experimentalist observing the computational system, and for which we shall define information flow. For now, we assume that we are interested in a single message.¹¹ We also assume that the message enters the computational system only at time $t = 0$, and at no later time instant.

We formally define the input nodes of the system as those nodes of \mathcal{G} , at time $t = 0$, whose transmissions statistically depend on the message M :

$$\mathcal{V}_{ip} := \{A_0 \in \mathcal{V}_0 : I(M; X(\mathbb{Q}(A_0))) > 0\}, \quad (3)$$

where $\mathbb{Q}(A_0)$ represents the set of edges leaving the node A_0 .

To remain consistent with Definition 3b, we define the computation performed by an input node $A_0 \in \mathcal{V}_{ip}$ as a function f_{A_0} that satisfies

$$f_{A_0}(M, W(A_0)) = X(\mathbb{Q}(A_0)), \quad (4)$$

and the computation performed by a non-input node at time $t = 0$, $A_0 \in \mathcal{V}_0 \setminus \mathcal{V}_{ip}$, as a function f_{A_0} that satisfies

$$f_{A_0}(W(A_0)) = X(\mathbb{Q}(A_0)). \quad (5)$$

As before, $W(A_0) \perp W(\mathcal{V}_0 \setminus \{A_0\}) \forall A_0 \in \mathcal{V}_0$ and $W(\mathcal{V}_0) \perp M$.

Remarks

1. Informally speaking, Definition 3 is designed to allow each node to generate a randomized function of its incoming transmissions for each of its outgoing transmissions.
2. The randomization at each node is explicitly captured by its intrinsic random variable $W(\cdot)$, and is assumed to be independent across all nodes of the system.
3. Furthermore, each node is allowed to send a different transmission on each of its outgoing edges.
4. Note that the condition imposed by Equation (2) introduces dependence between the random variables in the set $X(\mathcal{E})$.
5. For the most part, we will not be concerned with the precise form of the computation being performed by every node. We will only make use of information-theoretic measures applied to the message and to the random variables in the computational system.

Throughout the paper, we use the variables U, V, A, B, C and D to refer to nodes and E, P, Q, R and S to refer to edges. We use their script forms, e.g. \mathcal{R} , when referring to sets of nodes and edges, and primed script forms, e.g. \mathcal{R}' , when referring to subsets thereof. Once again, the notation we use is summarized in Table 1, and depicted in Figure 1 for convenience.

Having defined what we mean by the terms “computational system” and “message”, in the following sections we proceed to find a definition for “information flow” and identify properties that this definition satisfies in any computational system.

¹¹That is, we assume that the message is a single random variable or vector. It is possible to simultaneously examine the information flows of several (possibly dependent) messages, or of sub-messages within a single message. These cases are examined in Section 5.6.

TABLE 1
SUMMARY OF NOTATION

Variable(s)	Meaning
$\mathcal{G}^* = (\mathcal{V}^*, \mathcal{E}^*)$	The original complete directed graph, prior to time-unrolling
$\mathcal{G} = (\mathcal{V}, \mathcal{E})$	The time-unrolled graph making up the computational system
\mathcal{T}	The set of all time points, $\{0, 1, \dots, T\}$
\mathcal{V}	The set of all nodes in the computational system
\mathcal{V}_t	The subset of nodes at time t
V_t, A_t, B_t, C_t, D_t	A node in the graph at time t
V, A, B, C, D, E	A node in the original complete directed graph \mathcal{G}^* , or a node in the computational system at an unspecified time point
\mathcal{A}, \mathcal{B}	Some subset of nodes in \mathcal{V}
\mathcal{E}	The set of all edges in the computational system
\mathcal{E}_t	The set [†] of all edges at time t
\mathcal{E}'_t	Some subset [‡] of edges in \mathcal{E}_t
E_t, P_t, Q_t, R_t, S_t	An edge in the computational system at time t
E, P, Q, R, S	An edge in the original complete directed graph \mathcal{G}^* , or an edge in the computational system at an unspecified time point
$X(E_t)$	The random variable representing the transmission on the edge E_t
$X(\{E^{(1)}, E^{(2)}\})$	Short-hand notation for $\{X(E^{(1)}), X(E^{(2)})\}$ (refer Equation (1))
$\mathcal{P}(V_t)$	The set of all incoming edges of V_t ($= \mathcal{V}_{t-1} \times \{V_t\} \subseteq \mathcal{E}_{t-1}$)
$\mathcal{Q}(V_t)$	The set of all outgoing edges of V_t ($= \{V_t\} \times \mathcal{V}_{t+1} \subseteq \mathcal{E}_t$)
$W(V_t)$	The intrinsically generated random variable at the node V_t
M	The “message”, a random variable that enters the system at time $t = 0$, and whose information flow we seek to understand (refer Definition 3c)
\mathcal{V}_{ip}	The input nodes: the subset of nodes at time 0 whose outgoing transmissions depend on the message M (refer Definition 3c)
f_{V_t}	The function computed by the node V_t (refer Definition 3b)

[†]Script forms typically denote sets

[‡]Primed script forms typically denote subsets

3 Defining Information Flow

Before one can speak of *detecting* information flow in a network, it is first important to *define* what it is that we seek to detect.¹² In this section, we focus on arriving at a definition for information flow.

Our goal is to formalize how information about a message flows in a computational system. Ultimately, we expect to find the *path* that the message takes while being processed by the system. Towards this, we start by trying to formally define what it means for information about the message to flow on a given *edge*. This section concludes with a proposal for such a definition: one based on strict positivity of a conditional mutual information. But to provide the intuition behind this choice of definition, we start with several simpler candidate definitions, and show how they fail to satisfy an intuitive property using counterexamples.

After proposing a definition for information flow, in Section 4, we discuss the *properties* satisfied by our definition. Then, in Section 5, we specify how the transmissions of the computational system are observed, and describe how information flow might be *inferred* in a real computational system.

¹²In essence, “causal influence” measures such as Granger Causality and Directed Information, while intuitively quantifying transferred information, fail to lay down what *aspect of computation* they actually capture. This is, in part, a result of conflating the stages of defining a quantity we want to understand, and prescribing an estimator for it.

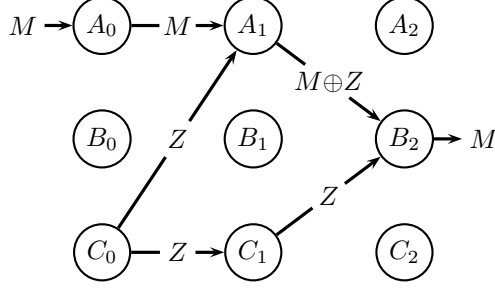


Figure 2: The computational system for Counterexample 1. We only depict edges relevant to the counterexample here. All other edges in the underlying complete directed graph are still present, but are not shown; their transmissions are assumed to be zero. Observe that no edge at time $t = 1$ has information flow as per Candidate Definition 1, yet the message reappears at time $t = 2$.

3.1 An intuitive property

To concretely define what it means for information about a message to flow on an edge, we need some way to assess competing candidate definitions and choose one among them. Towards this goal, we state a straightforward and intuitive property, which we would want any definition of information flow to satisfy.

Suppose that, at a given point in time, there is *no* flow of information about the message across *any* edge of a computational system. Note that this includes self-edges, so no node “carries” information about the message within its memory either. Then, we expect that information about the message has ceased to persist in the system, so the information flow about the message *must* be zero on all edges of the computational system, at all future points in time.

Property 1 (The Broken Telephone¹³). *Let \mathcal{C} be a computational system, and let $\mathcal{F}_M : \mathcal{E} \rightarrow \{0, 1\}$ be an indicator of the presence of information flow about M on an edge. That is, $\mathcal{F}_M(E) = 1$, if information about M flows on the edge $E \in \mathcal{E}$ and $\mathcal{F}_M(E) = 0$, otherwise. The Broken Telephone Property states that if, at some time $t \in \mathcal{T}$, we have*

$$\mathcal{F}_M(E_t) = 0 \quad \forall E_t \in \mathcal{E}_t, \tag{6}$$

then

$$\mathcal{F}_M(E_{t'}) = 0 \quad \forall E_{t'} \in \mathcal{E}_{t'}, \forall t' \in \mathcal{T}, t' > t. \tag{7}$$

3.2 Intuiting Information Flow through Counterexamples

We now propose four candidate definitions, beginning with the simplest. We then construct counterexamples to show how the first three candidate definitions do not satisfy Property 1.

Candidate Definition 1. *A simplistic and intuitive definition for information flow might simply stem from dependence. We say that information about the message M flows on an edge E_t if*

$$I(M; X(E_t)) > 0.$$

Counterexample 1. Consider the computational system depicted in Figure 2 (note that, in order to avoid unnecessary clutter, only edges with non-zero transmissions are shown in the figure). A_0 is the input node, which has the message $M \sim \text{Ber}(1/2)$ at time $t = 0$. The system’s goal is to communicate¹⁴ M to the

¹³https://en.wikipedia.org/wiki/Telephone_game

¹⁴This communication can be thought of as computing the identity function, and making the output available at the node B .

node B . It chooses the following strategy: at $t = 0$, A_0 “transmits” M to A_1 (i.e., node A stores M in its memory). C_0 independently generates a different random number, $W(C_0) = Z \sim \text{Ber}(1/2)$, $Z \perp M$, and sends this message to A_1 , while also storing it in memory until $t = 1$. A_1 then computes $M \oplus Z$ and passes the result to B_2 , while C_1 sends Z to B_2 . Here, the symbol “ \oplus ” stands for XOR, the exclusive-OR operator on two bits. B_2 is thus able to recover M by once again XOR-ing its inputs, $(M \oplus Z)$ and Z .

Note that the output of B_2 depends on M , even though none of its inputs individually depends on M . That is, $I(M; X((A_1, B_2))) = I(M; M \oplus Z) = 0$, and $I(M; X((C_1, B_2))) = I(M; Z) = 0$, so by Candidate Definition 1, information about the message flows on *no* edge at time $t = 1$. However, information about the message *does* flow out of node B_2 at time $t = 2$. This violates Property 1. Thus, mere *dependence* on the message cannot be a valid definition for flow of information on a single edge. \square

Communication strategies such as the one in Counterexample 1 frequently arise in cryptography [46], to prevent an eavesdropper from reading confidential information, and in network coding [12], for achieving the communication capacity of a network. Furthermore, a complex computational network may have smaller sub-networks with such topologies. For instance, we observe such a sub-network in the canonical example for network coding: the butterfly network [12, Fig. 7b] (this particular example is discussed in detail in Section 6.1). Optimal communication in such a network *requires* the use of such topologies, so Counterexample 1 is far from obscure. In fact, central to the idea of Counterexample 1 is a concept known as “synergy”, which is well-studied in the literature on Partial Information Decomposition [47–49] (see [50] for a recent review). This is discussed at length in Section 3.5. Even in neuroscience, the concept of synergy is recognized and well-understood [51–53], and some experimental evidence has appeared in the literature [54].

Counterexample 1 demonstrates that the information necessary to recover the message (or a function of it) is not necessarily transmitted through individual edges, but jointly across edges. So, we might instead seek to define the “smallest set of edges” along which information about the message flows, for every point in time. But if we ultimately wish to isolate *paths* along which information about the message flows, we require an understanding of which edges *specifically* the information flows upon. We therefore continue to think of information as flowing on individual edges.¹⁵

We can now update our naïve definition to counter the previous counterexample. We start by noting that in Counterexample 1, although the transmission on edge (A_1, B_2) is independent of M , it is not *conditionally* independent of M when given the transmission on (C_1, B_2) .

Candidate Definition 2. *We say that information about the message M flows on an edge $E_t \in \mathcal{E}_t$ if one of the following holds:*

1. $I(M; X(E_t)) > 0$, or
2. $\exists E'_t \in \mathcal{E}_t$ s.t. $I(M; X(E_t) | X(E'_t)) > 0$.

Counterexample 2. Consider a modified version of Counterexample 1, shown in Figure 3. Now, since there are *two* noise terms, no single extra edge may be conditioned upon to have non-zero information flow at time $t = 1$. So, Candidate Definition 2 also fails to satisfy Property 1. \square

It might seem that a possible rectification is to condition on *all* other edges at time t , but we can show that this also fails the test.

¹⁵It should be noted that the two views—information flowing on individual edges, versus sets of edges—are compatible with each other if we use Definition 5 (which will appear shortly) to describe information flow on a set of edges. This equivalence is elaborated upon in Section 3.4. Later, in Section 4.4, we attempt to refine our understanding of the aforementioned “smallest set of edges” along which information about the message flows.

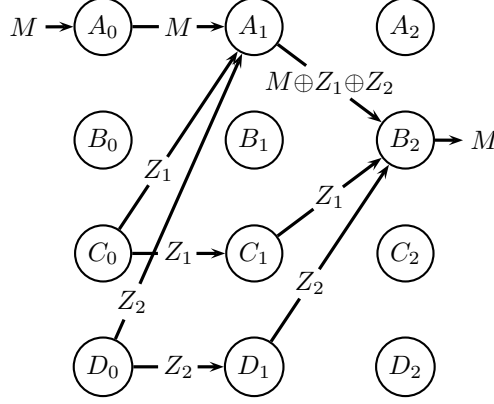


Figure 3: The computational system for Counterexample 2. Once again, observe that no edge at time $t = 1$ has information flow as per Candidate Definition 2, yet the message reappears at time $t = 2$. Note that only edges relevant to the counterexample are depicted in the figure. All other edges of the underlying complete directed graph are still present, and their transmissions are assumed to be zero.

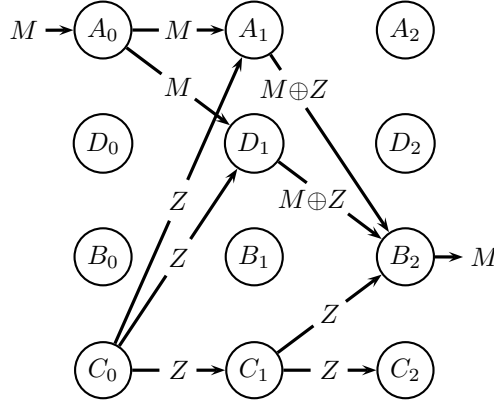


Figure 4: The computational system for Counterexample 3. Just as in the previous counterexamples, no edge at time $t = 1$ has information flow as per Candidate Definition 3, yet the message is reconstructed at time $t = 2$. Note that only edges relevant to the counterexample are depicted in the figure. All other edges of the underlying complete directed graph are still present, and their transmissions are assumed to be zero.

Candidate Definition 3. We say that information about the message M flows on an edge $E_t \in \mathfrak{E}_t$ if one of the following holds:

1. $I(M; X(E_t)) > 0$, or
2. $I(M; X(E_t) \mid X(\mathfrak{E}_t \setminus \{E_t\})) > 0$.

Counterexample 3. Consider the computational system shown in Figure 4. Once again, we have an input node A_0 which possesses the message at time $t = 0$, and wishes to send this message to node B . It does so by mixing M with an independent random variable Z generated at C_0 , so that the scenario described in Counterexample 1 still holds. But additionally, A communicates to B along a redundant path, through D_1 . Now, if E is any incoming edge of B_2 , it is still true that $I(M; X(E)) = 0$. So none of the inputs of B_2 individually depends on M , thus eliminating the first condition in Candidate Definition 3. Furthermore, checking each incoming edge of B_2 reveals that the second condition also fails to hold. If we take $E_1 = (A_1, B_2)$, we get

$$I(M; X(E_1) \mid X(\mathfrak{E}_1 \setminus \{E_1\})) = I(M; M \oplus Z \mid M \oplus Z, Z) = 0. \quad (8)$$

The same holds true when $E_1 = (D_1, B_2)$ since the transmissions on both edges are identical by construction. Likewise, if we take $E_1 = (C_1, B_2)$, we have

$$I(M; X(E_1) \mid X(\mathcal{E}_1 \setminus \{E_1\})) = I(M; Z \mid M \oplus Z, Z) = 0, \quad (9)$$

with the same holding true when $E_1 = (C_1, C_2)$. Therefore, no edge at time $t = 1$ has any information flow about the message M , as per Candidate Definition 3. Nevertheless, B_2 is able to recover the message at time $t = 2$, proving that Property 1 fails to hold for Candidate Definition 3. \square

3.3 Information Flow on a Single Edge

The counterexamples presented in the previous section motivate a new definition for when information about the message can be said to flow on a given edge. Neither *dependence* of M on the transmission of an edge, nor conditional dependence given *one* or *all* other edges, satisfy Property 1.

However, in all these counterexamples, given an edge E_t upon which we expect to have non-zero information flow, we observe: there is at least one subset of edges $\mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\}$, such that when given $X(\mathcal{E}'_t)$, $X(E_t)$ is conditionally dependent¹⁶ on M . In Counterexample 1, the edge (A_1, B_2) , carrying $M \oplus Z$, is conditionally dependent on M , given $X((C_1, B_2)) = Z$. In Counterexample 2, $X((A_1, B_2)) = M \oplus Z_1 \oplus Z_2$ is conditionally dependent on M , given $\{X((C_1, B_2)), X((D_1, B_2))\} = \{Z_1, Z_2\}$. And finally, in Counterexample 3, $X((A_1, B_2)) = M \oplus Z$ is conditionally dependent on M , given $X((C_1, B_2)) = Z$; note that we do *not* condition on $X((D_1, B_2)) = M \oplus Z$. Thus, conditioning on a subset of the other edges' transmissions creates dependence between M and the transmission on an edge of interest.

We will shortly prove that Property 1 holds when information flow is defined as below, so we directly state it as a definition, skipping its candidacy status.

Definition 4 (*M-information Flow on a Single Edge*). *We say that information about the message M flows on an edge $E_t \in \mathcal{E}_t$ if*

$$\exists \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\} \quad \text{s.t.} \quad I(M; X(E_t) \mid X(\mathcal{E}'_t)) > 0. \quad (10)$$

Henceforth, we refer to “information flow about the message M ” as M -information flow, and use the phrase “the edge E_t has M -information flow” or “the edge E_t carries M -information flow” to mean that information about M flows on E_t per this definition.

Note that if $I(M; X(E_t) \mid X(\mathcal{E}'_t)) > 0$, then $I(M; X(\{E_t\} \cup \mathcal{E}'_t)) > 0$. In other words, *there exists* a set of edges that includes E_t , whose transmissions depend on M . This is why it is important to condition on all possible subsets of \mathcal{E}_t . It is not immediately clear, however, whether *every* edge in $\{E_t\} \cup \mathcal{E}'_t$ has M -information flow. We return to this point in Section 4.4.

Also, this definition implies that certain edges, such as (C_1, B_2) in Counterexample 1, may have M -information flow, which may seem counter-intuitive. This is discussed further and justified in Section 4.2.

3.4 Information Flow on a Set of Edges

The definition of M -information flow for a single edge naturally generalizes to one for a set of edges, at a given time.

¹⁶Equivalently, we could say that there exists at least one subset of edges $\mathcal{E}'_t \subseteq \mathcal{E}_t$, without explicitly excluding E_t , since $I(M; X(E_t) \mid X(E_t), X(\mathcal{E}'_t)) = 0$.

Definition 5 (*M*-information Flow on a Set of Edges). We say that information about the message M flows on a set of edges $\mathcal{E}'_t \subseteq \mathcal{E}_t$ if

$$\exists \mathcal{R}'_t \subseteq \mathcal{E}_t \quad \text{s.t.} \quad I(M; X(\mathcal{E}'_t) | X(\mathcal{R}'_t)) > 0. \quad (11)$$

The definition of *M*-information flow on a set of edges is nearly identical to its single-edge counterpart. Indeed, they are closely related, as the following proposition shows.

Proposition 1. A set $\mathcal{E}'_t \subseteq \mathcal{E}_t$ has *M*-information flow if and only if there exists an edge $E'_t \in \mathcal{E}'_t$ that has *M*-information flow.

A proof of this proposition can be found in Appendix A.

It should be noted that although the counterexamples in this section all employed computational systems which recovered the message M at a new node at a later time, a computational system will in general compute some function of the message. For instance, see the example in Section 6.2.

3.5 The Connection with Synergistic Information

This section connects our definition of *M*-information flow with recent developments on a subject known as “Partial Information Decomposition” (PID). Our definition is closely related to the concept of “Synergistic Information” that appears in this field. This section exists only for the purpose of providing a deeper intuition for our definition of *M*-information flow, and does not affect the rest of the paper in any significant way. We have attempted to explain this intuition in a way that is accessible to readers unfamiliar the PID literature. However, readers may feel free to skip this section, if desired.

At its core, Counterexample 1 relies on a concept known as “synergy”, which is described explicitly in the literature on Partial Information Decomposition (PID) [47–49] (see [50] for a recent review, and Appendix C for a brief introduction). Essentially, this body of literature seeks to decompose the mutual information that two or more variables share about a message, $I(M; (Y_1, Y_2, \dots))$, into several individually meaningful, non-negative components. In particular, when discussing the bivariate case—i.e., the case of two variables, $I(M; (Y_1, Y_2))$ —it is understood what the terms in this decomposition should be: (i) information about the message that each variable carries *uniquely*, and which cannot be inferred from the other; (ii) information about the message that the variables share *redundantly*, and which can be extracted from either; (iii) and information about the message that the variables convey *synergistically*, which is revealed only when *both* variables are taken together, and cannot be inferred from either variable *individually*. Counterexample 1 is the canonical example for synergy, and is known simply as the “XOR” example in the PID literature. While $M \oplus Z$ and Z are *individually* independent of M , when taken *together*, $I(M; (M \oplus Z, Z)) = H(M)$. This suggests that $M \oplus Z$ and M have no unique or shared information about M , but convey information synergistically.

While the field has not yet arrived at a consensus on the most appropriate definitions for unique, redundant and synergistic information [50], it is well-understood what properties these quantities must satisfy, at least in the bivariate case (see Appendix C, specifically, Equations (94), (95) and (97)). Therefore, even without formal definitions, we can rely on the intuition provided by these properties to understand the implications of PID for *M*-information flow. If a particular edge’s transmission contains unique or redundant information about the message (with respect to some other subset of edges at that point in time), then that information will manifest itself in the form of strictly positive mutual information. However, in the absence of positive mutual information between the message and the transmission on a given edge, we need to consider whether said transmission synergistically interacts with another subset of transmissions at that point in time, as this could potentially create dependence with the message through the kind of “recombination” described in

Counterexample 1. We then need to decide whether such synergistic interactions ought to be considered to constitute information flow. As we show below, our definition of M -information flow *does* consider instances of purely synergistic information to constitute information flow.

Indeed, it is possible to formulate a definition for information flow based on synergy, which is completely equivalent to Definition 4. The definition below makes use of the PID preliminaries given in Appendix C.

Definition 6 (M -synergistic information flow). *We say that an edge E_t has M -synergistic information flow if at least one of the following holds:*

1. $I(M; X(E_t)) > 0$, or
2. $\exists \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\}$ s.t. $CI(M : X(E_t); X(\mathcal{E}'_t)) > 0$,

where $CI(M : X; Y)$ represents the synergistic information between X and Y about M .

Proposition 2 (Equivalence of Information Flow Definitions). *An edge E_t has M -information flow if and only if it has M -synergistic information flow. Furthermore, suppose E_t is an edge which satisfies $I(M; X(E_t)) = 0$. Then,*

$$I(M; X(E_t) | X(\mathcal{E}'_t)) > 0 \tag{12}$$

for some set $\mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\}$, if and only if

$$CI(M : X(E_t); X(\mathcal{E}'_t)) > 0. \tag{13}$$

That is, the set \mathcal{E}'_t upon whose transmissions we need to condition is the same as the one responsible for providing synergy in the alternate definition.

A proof of this proposition is given in Appendix C.

We should also mention here that it may be possible to leverage the more recent definitions of synergy to supply an intuitive measure of the *volume* of information flow; we discuss this in Section 7.5.

4 Properties of Information Flow

Having defined what it means for information about a message to flow on an edge, we demonstrate that Definition 4 satisfies several intuitively desirable properties, including Property 1.

4.1 The Broken Telephone Property

Theorem 3. *M -information flow, as given by Definition 4, satisfies Property 1.*

Before we prove this theorem, we prove a simpler lemma which directly falls out of Definition 4 and the properties of mutual information.

Lemma 4. *There is no edge in \mathcal{E}_t that carries M -information flow if, and only if, $X(\mathcal{E}_t)$ is independent of M . In other words,*

$$I(M; X(E_t) | X(\mathcal{E}'_t)) = 0 \quad \forall E_t \in \mathcal{E}_t, \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\} \tag{14}$$

if and only if

$$I(M; X(\mathcal{E}_t)) = 0. \tag{15}$$

Equivalently, we can state the opposite: $X(\mathcal{E}_t)$ depends on M if and only if at least one edge in \mathcal{E}_t carries M -information flow.

Proof. (\Rightarrow) Suppose that the condition in (14) holds. Let $\mathfrak{E}_t = \{E_t^{(1)}, E_t^{(2)}, \dots, E_t^{(N^2)}\}$ be any ordering of the edges in \mathfrak{E}_t . Then,

$$I(M; X(\mathfrak{E}_t)) \stackrel{(a)}{=} I(M; X(E_t^{(1)})) + I(M; X(E_t^{(2)}) | X(E_t^{(1)})) \quad (16)$$

$$+ I(M; X(E_t^{(3)}) | X(E_t^{(1)}), X(E_t^{(2)})) + \dots$$

$$= \sum_{i=1}^{N^2} I\left(M; X(E_t^{(i)}) \mid \bigcup_{j=1}^{i-1} \{X(E_t^{(j)})\}\right) \quad (17)$$

$$\stackrel{(b)}{=} \sum_{i=1}^{N^2} I\left(M; X(E_t^{(i)}) \mid X\left(\bigcup_{j=1}^{i-1} \{E_t^{(j)}\}\right)\right) \stackrel{(c)}{=} 0, \quad (18)$$

where (a) follows from the chain-rule of mutual information [55, Ch. 2], (b) is simply the application of Equation (1), and (c) follows from the fact that each term in the summation is zero, by (14). This proves the forward implication.

(\Leftarrow) Next, suppose $I(M; X(\mathfrak{E}_t)) = 0$. Let E_t be any edge in \mathfrak{E}_t and let \mathfrak{E}'_t be any subset of $\mathfrak{E}_t \setminus \{E_t\}$. Also, let $\mathfrak{E}''_t = \mathfrak{E}_t \setminus (\mathfrak{E}'_t \cup \{E_t\})$. Then,

$$0 = I(M; X(\mathfrak{E}_t)) \quad (19)$$

$$= I(M; X(\mathfrak{E}'_t)) + I(M; X(E_t) | X(\mathfrak{E}'_t)) + I(M; X(\mathfrak{E}''_t) | X(\mathfrak{E}'_t), X(E_t)) \quad (20)$$

by the chain rule. Since (conditional) mutual information is always non-negative [55, Ch. 2], all three terms on the right hand side must be zero. So in particular,

$$I(M; X(E_t) | X(\mathfrak{E}'_t)) = 0. \quad (21)$$

Since E_t and \mathfrak{E}'_t are arbitrary, this proves the converse. \square

Proof of Theorem 3. We need to prove that M -information flow, as given by Definition 4, satisfies Property 1. Explicitly stated, we need to show that if every edge at some time t has zero M -information flow, then every edge at all future times $t' > t$ must also have zero M -information flow. So suppose that, at time t , for every $E_t \in \mathfrak{E}_t$ we have

$$I(M; X(E_t) | X(\mathfrak{E}'_t)) = 0 \quad \forall \mathfrak{E}'_t \subseteq \mathfrak{E}_t \setminus \{E_t\}. \quad (22)$$

By Lemma 4, this implies that

$$I(M; X(\mathfrak{E}_t)) = 0. \quad (23)$$

Now, consider the first future time instant, $t' = t + 1$. For every node $A_{t+1} \in \mathcal{V}_{t+1}$, the definition of computation at a node (Definition 3b) states that

$$X(\mathbb{Q}(A_{t+1})) = f_{A_{t+1}}(X(\mathcal{P}(A_{t+1})), W(A_{t+1})), \quad (24)$$

where the reader may recall, $\mathcal{P}(A_{t+1})$ and $\mathbb{Q}(A_{t+1})$ are the edges entering and leaving A_{t+1} respectively. We can collect the individual functions $f_{A_{t+1}}$ across all nodes in \mathcal{V}_{t+1} into a single joint function $f_{\mathcal{V}_{t+1}}$, as described in Definition 3b, to obtain

$$X(\mathfrak{E}_{t+1}) = f_{\mathcal{V}_{t+1}}(X(\mathfrak{E}_t), W(\mathcal{V}_{t+1})). \quad (25)$$

Therefore,

$$0 \leq I(M; X(\mathfrak{E}_{t+1})) = I(M; f_{\mathcal{V}_{t+1}}(X(\mathfrak{E}_t), W(\mathcal{V}_{t+1}))) \quad (26)$$

$$\stackrel{(a)}{\leq} I(M; X(\mathfrak{E}_t), W(\mathcal{V}_{t+1})) \quad (27)$$

$$= I(M; X(\mathfrak{E}_t)) + I(M; W(\mathcal{V}_{t+1}) | X(\mathfrak{E}_t)) \quad (28)$$

$$\stackrel{(b)}{=} I(M; X(\mathfrak{E}_t)) \stackrel{(c)}{=} 0, \quad (29)$$

where (a) follows from the Data Processing Inequality [55, Ch. 2], (b) follows from the fact that $W(\mathcal{V}_{t+1}) \perp \{M, X(\mathcal{E}_t)\}$, and (c) follows from (23). Once again, by non-negativity of mutual information we must have that $I(M; X(\mathcal{E}_{t+1})) = 0$. Applying Lemma 4 once again, we find that for $t' = t + 1$,

$$I(M; X(E_{t'}) | X(\mathcal{E}'_{t'})) = 0 \quad \forall E_{t'} \in \mathcal{E}_{t'}, \mathcal{E}'_{t'} \subseteq \mathcal{E}_{t'} \setminus \{E_{t'}\} \quad (30)$$

We have shown that (22) implies (30), so induction on t' yields that (30) holds for all future times $t' > t$, completing the proof. \square

4.2 The Existence of Orphans

Definition 4 also has a very non-intuitive property: an edge leading out of a node may have M -information flow, even though *no* edge leading *into* that node has M -information flow.

Definition 7 (M -information Orphan). *In a computational system \mathcal{C} , a node V_t is said to be an M -information orphan if $\mathcal{Q}(V_t)$ has M -information flow (as per Definition 5), but $\mathcal{P}(V_t)$ has no M -information flow.*

Property 2. *M -information orphans may exist in a computational system.*

Proof. Consider the computational system in Figure 2 from Counterexample 1. The node C_1 is an M -information orphan, since the edge (C_1, B_2) carries M -information flow, whereas none of its incoming edges carries M -information flow. \square

The existence of M -information orphans, along with the presence of M -information flow on (C_1, B_2) in Counterexample 1, may not be expected, since Z was never computed from M . Indeed, M -information flow appears to emerge from “nowhere” at the node C_1 , leaving nodes such as C_1 orphaned in a view of the graph that contains only edges having M -information flow, and hence the name. But closer inspection reveals that in this example, the transmissions arriving at B_2 from A_1 and C_1 , i.e. $M \oplus Z$ and Z , are *statistically identical*: they are both individually independent of M , but when XOR’ed, are fully dependent on M . In other words, any *purely observational* measure¹⁷ defined on the transmissions at time t that assigns M -information flow to $M \oplus Z$, must also assign M -information flow to Z .

Note that, just as M -information flow can originate at an M -information orphan, M -information flow may also terminate at a node—either by simple omission, or as a result of some computation (see Section 6 for such instances). Likewise, multiple outgoing edges of a given node may transmit redundant copies of the same information. Ultimately, we see that there is no “law of conservation” for M -information flow. In this sense, “information flow” is not a typical kind of “flow” that is defined on graphs (see, for example, [56, Sec. 26.1]), and well-known results such as the Max-flow Min-cut Theorem [56, Thm. 26.6] do not apply as-is to M -information flow.

It is worthwhile to note at this point that the existence of M -information orphans such as C_1 in Counterexample 1 is not inconsistent with the Data Processing Inequality [55, Ch. 2]. In fact, a clear example of the Data Processing Inequality in play is seen at the network-level, wherein $M - X(\mathcal{E}_t) - X(\mathcal{E}_{t+1})$ form a Markov Chain for any time $0 \leq t < T$, and so the information content about M present collectively in all transmissions at time $t + 1$ *must* be no more than that present at time t . We call this Global Markovity, and state it formally for completeness.

Corollary 5 (Global Markovity). *At any given time t , the following Markov Chain holds: $M - X(\mathcal{E}_t) - X(\mathcal{E}_{t+1})$.*

¹⁷i.e., a functional of the joint distribution of $X(\mathcal{E}_t)$

In fact, this Markov condition must hold for every *subset* of nodes, not just for the entire set of nodes, so it is subsumed by the following proposition.

Proposition 6 (Local Markovity). *For any given subset of nodes $\mathcal{V}'_t \subseteq \mathcal{V}_t$, the following Markov Chain holds: $M - X(\mathcal{P}(\mathcal{V}'_t)) - X(\mathcal{Q}(\mathcal{V}'_t))$.*

Proof. Since $X(\mathcal{Q}(\mathcal{V}'_t)) = f_{\mathcal{V}'_t}(X(\mathcal{P}(\mathcal{V}'_t)), W(\mathcal{V}'_t))$ by Definition 3b, the tuple $(X(\mathcal{P}(\mathcal{V}'_t)), X(\mathcal{Q}(\mathcal{V}'_t)))$ is also a function of $X(\mathcal{P}(\mathcal{V}'_t))$ and $X(W(\mathcal{V}'_t))$. Hence, the following Markov chain holds:

$$M - (X(\mathcal{P}(\mathcal{V}'_t)), W(\mathcal{V}'_t)) - (X(\mathcal{P}(\mathcal{V}'_t)), X(\mathcal{Q}(\mathcal{V}'_t))).$$

By the Data Processing Inequality, this implies that

$$I(M; X(\mathcal{Q}(\mathcal{V}'_t)), X(\mathcal{P}(\mathcal{V}'_t))) \leq I(M; X(\mathcal{P}(\mathcal{V}'_t)), W(\mathcal{V}'_t)) \quad (31)$$

$$\stackrel{(a)}{=} I(M; X(\mathcal{P}(\mathcal{V}'_t))) + I(M; W(\mathcal{V}'_t) \mid X(\mathcal{P}(\mathcal{V}'_t))) \quad (32)$$

$$\stackrel{(b)}{=} I(M; X(\mathcal{P}(\mathcal{V}'_t))) + I(W(\mathcal{V}'_t); M, X(\mathcal{P}(\mathcal{V}'_t))) - I(W(\mathcal{V}'_t); X(\mathcal{P}(\mathcal{V}'_t))) \quad (33)$$

$$\stackrel{(c)}{=} I(M; X(\mathcal{P}(\mathcal{V}'_t))) + 0 - 0, \quad (34)$$

where in (a) and (b), we have used the chain rule of mutual information in two different ways, and in (c) we have used the fact that $W(\mathcal{V}'_t) \perp \{M, X(\mathcal{P}(\mathcal{V}'_t))\}$. Therefore,

$$I(M; X(\mathcal{Q}(\mathcal{V}'_t)) \mid X(\mathcal{P}(\mathcal{V}'_t))) = 0, \quad (35)$$

which implies the Markov chain in Proposition 6. \square

Since the above also holds for $\mathcal{V}'_t = \mathcal{V}_t$, wherein $\mathcal{Q}(\mathcal{V}_t) = \mathcal{E}_t$, Proposition 6 implies Corollary 5.

Given that these Markov conditions arise directly from the way we have defined the computational system, specifically Definition 3b, they may not be very surprising (indeed, they may be considered *properties* of the computational system model itself). However, it is worth noting that Proposition 6 holds *even at an M-information orphan*. Thus, M -information orphans do not “create” information about M , as we would rightly expect, given the Data Processing Inequality.

4.3 The Existence of Information Paths

We now show that if the outgoing transmissions of any given node depend on the message, then we can find a path leading to that node from one or more input nodes, along which M -information flows. Before we demonstrate this property, we formally define what we mean by the terms “path” and “cut”.

Definition 8 (Path). *In any computational system \mathcal{C} , suppose \mathcal{A} and \mathcal{B} are two disjoint sets of nodes in \mathcal{V} . Then, a path from \mathcal{A} to \mathcal{B} is any ordered set of nodes $\{V^{(0)}, V^{(1)}, \dots, V^{(L)}\}$ that satisfies (i) $V^{(0)} \in \mathcal{A}$; (ii) $V^{(L)} \in \mathcal{B}$; and (iii) $(V^{(i-1)}, V^{(i)}) \in \mathcal{E}$ for every $1 \leq i \leq L$, where L is a positive integer indicating the length of the path. We refer to the set $\{(V^{(i-1)}, V^{(i)})\}_{i=1}^L$ as the edges of the path.*

Definition 9 (M -Information Path). *Continuing from Definition 8, we define an M -information path from \mathcal{A} to \mathcal{B} as any path from \mathcal{A} to \mathcal{B} , each of whose edges carries M -information flow. That is, if $(V^{(i-1)}, V^{(i)}) = E_{t_i} \in \mathcal{E}_{t_i}$ for some $t_i \in \mathcal{T}$, then for every $1 \leq i \leq L$,*

$$\exists \mathcal{E}'_{t_i} \subseteq \mathcal{E}_{t_i} \quad \text{s.t.} \quad I(M; X(E_{t_i}) \mid X(\mathcal{E}'_{t_i})) > 0. \quad (36)$$

Definition 10 (Cut). In any computational system \mathcal{C} , suppose \mathcal{A} and \mathcal{B} are two disjoint sets of nodes in \mathcal{V} . Then, a cut separating \mathcal{A} and \mathcal{B} is any pair of sets $(\mathcal{V}^{\text{src}}, \mathcal{V}^{\text{sink}})$, such that (i) $\mathcal{V}^{\text{src}} \cup \mathcal{V}^{\text{sink}} = \mathcal{V}$; (ii) $\mathcal{V}^{\text{src}} \cap \mathcal{V}^{\text{sink}} = \emptyset$; (iii) $\mathcal{A} \subseteq \mathcal{V}^{\text{src}}$; and (iv) $\mathcal{B} \subseteq \mathcal{V}^{\text{sink}}$. We refer to the set of edges going from \mathcal{V}^{src} to $\mathcal{V}^{\text{sink}}$, i.e. $\mathcal{E} \cap (\mathcal{V}^{\text{src}} \times \mathcal{V}^{\text{sink}})$, as the edges in the cut set¹⁸.

Definition 11 (Zero- M -information Cut). Continuing from Definition 10, we say that a cut $(\mathcal{V}^{\text{src}}, \mathcal{V}^{\text{sink}})$ is a zero- M -information cut if every edge in its cut set has zero M -information flow. That is, for every $E_t \in \mathcal{E} \cap (\mathcal{V}^{\text{src}} \times \mathcal{V}^{\text{sink}})$,

$$I(M; X(E_t) \mid X(\mathcal{E}'_t)) = 0 \quad \forall \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\}. \quad (37)$$

Remark In Definition 11, we require that Equation (37) hold for every edge E_t in $\mathcal{E} \cap (\mathcal{V}^{\text{src}} \times \mathcal{V}^{\text{sink}})$. However, the edges in this set may belong to several different time points, since the cut is not restricted to any particular time (e.g., see Figure 5). The time t used in Equation (37), therefore, is determined by the time of the edge E_t , and varies for each E_t that we check in $\mathcal{E} \cap (\mathcal{V}^{\text{src}} \times \mathcal{V}^{\text{sink}})$.

Property 3 (Existence of an Information Path). In any computational system \mathcal{C} , suppose that at some time $t_{op} \in \mathcal{T}$, there is an “output node” $V_{op} \in \mathcal{V}$ whose outgoing edges $\mathcal{Q}(V_{op})$ satisfy $I(M; X(\mathcal{Q}(V_{op}))) > 0$. Then, there must exist an M -information path from the input nodes \mathcal{V}_{ip} to V_{op} .

Theorem 7. Definition 4 satisfies Property 3.

While the theorem seems obvious on the surface, the proof is in fact non-trivial because of the nature of our definition of M -information flow. Due to Property 2, M -information flowing *out of* a node does *not* imply that M -information must flow *into* that node. Therefore, a straightforward application of the Data Processing Inequality at every node fails to prove the theorem, and we must resort to a more rigorous cut-set-based approach.

Proof outline. We shall prove the contrapositive of the theorem, i.e., we will show that if there exists no M -information path from \mathcal{V}_{ip} to V_{op} , then the outgoing transmissions of V_{op} are independent of M . We first connect the absence of any M -information path with the presence of a zero- M -information cut. This is achieved in Lemma 8, which we present before the proof of Theorem 7.

The proof itself proceeds by induction over time. We divide the proof into two steps: initialization and continuation. Starting with the first nodes that come after the cut (temporally) in the initialization step, we systematically show that all nodes to the right of the cut have outgoing transmissions that are independent of the message M through induction. In this proof outline, we show these steps intuitively using Figure 5, where the dashed black line denotes the cut.

Initialization. Here, node C_1 is the first node to the right of the cut, and all of its incoming edges must come from across the cut (depicted by lines in red). Because the cut is a zero- M -information cut, none of its incoming transmissions have M -information flow. Furthermore, the intrinsically generated random variable $W(C_1)$ is independent of M . Using these two facts along with the Data Processing Inequality, we can show that the transmissions on C_1 ’s outgoing edges, $X(\mathcal{Q}(C_1))$, are also independent of M .

Continuation. At the second time instant to the right of the cut, nodes B_2 and C_2 receive their incoming transmissions from either C_1 (shown in orange) or from across the cut (shown in blue). Once again, the transmissions coming from across the cut can have no information flow, and we have shown that the transmissions coming from C_1 are independent of M . Also, $W(B_2)$ and $W(C_2)$ are independent of M and all incoming transmissions. This suffices to show that the outgoing transmissions of B_2 and C_2 , $X(\mathcal{Q}(B_2) \cup \mathcal{Q}(C_2))$, are

¹⁸Note that it is not necessary for us to assume that, individually, \mathcal{V}^{src} and $\mathcal{V}^{\text{sink}}$ are *connected* sets of nodes. For instance, there may be an isolated subset of $\mathcal{V}^{\text{sink}}$, surrounded only by nodes in \mathcal{V}^{src} . Our theorems and proofs remain unaffected, even in such a scenario.

independent of M . Applying this argument repeatedly over time shows that the transmissions of all nodes to the right of the cut are independent of M .

Therefore, if there is a node V_{op} whose outputs depend on M , we can be assured that there exists no zero- M -information cut separating \mathcal{V}_{ip} from V_{op} . Therefore, by Lemma 8, there exists an M -information path from \mathcal{V}_{ip} to V_{op} . \square

A few nuances are omitted in this outline, such as how the definition of \mathcal{V}_{ip} plays a role precisely. These subtleties are better elucidated in the full proof.

Before proceeding to the formal proof of Theorem 7, we first state and prove the lemma we alluded to earlier, which shows how the absence of an M -information path implies the presence of a zero- M -information cut, and vice versa.

Lemma 8. *Let \mathcal{A} and \mathcal{B} be two disjoint sets of nodes in the computational system \mathcal{C} . There exists no M -information path from \mathcal{A} to \mathcal{B} if and only if there is a zero- M -information cut separating \mathcal{A} and \mathcal{B} .*

Proof. (\Rightarrow) Suppose there exists no M -information path from \mathcal{A} to \mathcal{B} . Consider the set of all nodes to which there exists at least one M -information path from \mathcal{A} . Let \mathcal{V}^{src} be the collection of all such nodes, along with the nodes in \mathcal{A} , i.e.,

$$\mathcal{V}^{\text{src}} := \mathcal{A} \cup \{V_t \in \mathcal{V} : \exists \text{ an } M\text{-information path from } \mathcal{A} \text{ to } V_t\}. \quad (38)$$

Let $\mathcal{V}^{\text{sink}} = \mathcal{V} \setminus \mathcal{V}^{\text{src}}$, so that $\mathcal{V}^{\text{sink}}$ consists of nodes to which there is no M -information path from \mathcal{A} . Then, we must have $\mathcal{B} \subseteq \mathcal{V}^{\text{sink}}$, since it is known that there are no M -information paths from \mathcal{A} to \mathcal{B} . Therefore, $(\mathcal{V}^{\text{src}}, \mathcal{V}^{\text{sink}})$ is a cut that separates \mathcal{A} and \mathcal{B} , such that no edge in the cut set has M -information flow. In other words, by Definition 11, this is a zero- M -information cut separating \mathcal{A} and \mathcal{B} .

(\Leftarrow) Next, suppose that there *is* an M -information path $\{V^{(i)}\}_{i=0}^L$ from \mathcal{A} to \mathcal{B} . Then, we claim that there can exist no zero- M -information cut separating \mathcal{A} and \mathcal{B} . Let $(\mathcal{V}^{\text{src}}, \mathcal{V}^{\text{sink}})$ be any cut separating \mathcal{A} and \mathcal{B} . By Definition 8, we must have $V^{(0)} \in \mathcal{V}^{\text{src}}$ and $V^{(L)} \in \mathcal{V}^{\text{sink}}$. So, there must be at least one edge going from \mathcal{V}^{src} to $\mathcal{V}^{\text{sink}}$ which lies on the path. This implies that at least one edge in the cut set carries M -information flow. Since the conditions of Definition 11 are not satisfied, this cut is *not* a zero- M -information cut. Since this is true for every cut separating \mathcal{A} and \mathcal{B} , the claim holds. \square

Proof of Theorem 7. As mentioned in the proof outline, we prove the contrapositive of the theorem. Suppose there exists no M -information path from the input nodes \mathcal{V}_{ip} to V_{op} . Then, by Lemma 8, there exists a zero- M -information cut¹⁹ separating \mathcal{V}_{ip} and V_{op} . We use this to prove that the transmissions of V_{op} are independent of M .

Setup. Let the cut separating \mathcal{V}_{ip} and V_{op} be given by $(\mathcal{V}^{\text{src}}, \mathcal{V}^{\text{sink}})$, so that $\mathcal{V}_{\text{ip}} \subseteq \mathcal{V}^{\text{src}}$ and $V_{\text{op}} \in \mathcal{V}^{\text{sink}}$. Then, the cut divides \mathcal{E} into the following sets: $\mathcal{E}^{\text{src}} = \mathcal{E} \cap (\mathcal{V}^{\text{src}} \times \mathcal{V}^{\text{src}})$, the edges between the nodes in \mathcal{V}^{src} ; $\mathcal{E}^{\text{sink}} = \mathcal{E} \cap (\mathcal{V}^{\text{sink}} \times \mathcal{V}^{\text{sink}})$, the edges between nodes in $\mathcal{V}^{\text{sink}}$; and $\mathcal{E}^{\text{cut}} = \mathcal{E} \cap (\mathcal{V}^{\text{src}} \times \mathcal{V}^{\text{sink}})$, the edges going from \mathcal{V}^{src} to $\mathcal{V}^{\text{sink}}$ (the edges going from $\mathcal{V}^{\text{sink}}$ to \mathcal{V}^{src} will not be relevant to our discussion). From the previous paragraph, Lemma 8 implies that $(\mathcal{V}^{\text{src}}, \mathcal{V}^{\text{sink}})$ is a zero- M -information cut, so by Definition 11, we have that for all $E_t \in \mathcal{E}^{\text{cut}}$,

$$I(M; X(E_t) \mid X(\mathcal{E}'_t)) = 0 \quad \forall \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\}. \quad (39)$$

Note that the edges in \mathcal{E}^{cut} may belong to different time instants. In particular, the time instant t in the equation above corresponds to the time of the edge E_t , whose flow is in question.²⁰

¹⁹Note that, in general, this cut may be arbitrarily complex, spanning several nodes and multiple time instants.

²⁰In fact, this is one of the central factors that prevents us from recursively applying the Data Processing Inequality at every node, leading from \mathcal{V}_{ip} to V_{op} .

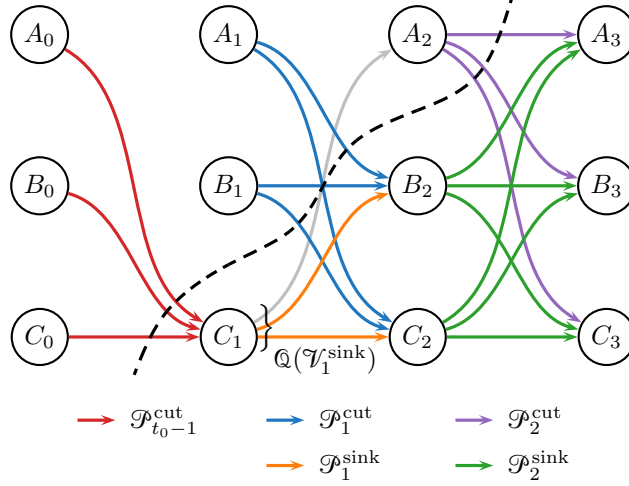


Figure 5: A generic computational system used in the proof outline and to explain certain steps in the proof of Theorem 7. For the purposes of the proof outline, it suffices to note that the black dashed line denotes the cut. All variable names can be ignored at this point of time.

For the purposes of the formal proof, note that in this figure, \mathcal{E}^{cut} is essentially the union of the red, blue and purple edges, while $\mathcal{E}^{\text{sink}}$ is the union of the orange and green edges. From this, it is evident that $\mathcal{P}(\mathcal{V}_t^{\text{sink}}) = \mathcal{P}_{t-1}^{\text{cut}} \cup \mathcal{P}_{t-1}^{\text{sink}}$ for any time t , i.e., the incoming edges of $\mathcal{V}^{\text{sink}}$ at time t must either come from nodes in $\mathcal{V}^{\text{sink}}$ or from nodes across the cut. Secondly, it should be clear that $\mathcal{P}_{t-1}^{\text{sink}} = \mathcal{Q}(\mathcal{V}_{t-1}^{\text{sink}}) \cap \mathcal{E}^{\text{sink}}$, i.e., the *incoming* edges of $\mathcal{V}_t^{\text{sink}}$ that originate from nodes in $\mathcal{V}^{\text{sink}}$ are simply the *outgoing* edges of $\mathcal{V}_{t-1}^{\text{sink}}$ which terminate at nodes in $\mathcal{V}^{\text{sink}}$. This is seen best at time $t = 1$ in the graph above, where the orange and grey lines together represent $\mathcal{Q}(\mathcal{V}_1^{\text{sink}})$, the orange and green edges together make up $\mathcal{E}^{\text{sink}}$, and $\mathcal{P}_1^{\text{sink}}$ is given by the orange edges, which is the intersection of the two sets.

Order the nodes in $\mathcal{V}^{\text{sink}}$ by time, and let $\mathcal{V}_t^{\text{sink}}$ be the subset of nodes in $\mathcal{V}^{\text{sink}}$ at time t . Let $\mathcal{P}(\mathcal{V}_t^{\text{sink}})$ and $\mathcal{Q}(\mathcal{V}_t^{\text{sink}})$ respectively be the sets of edges *collectively* entering and leaving all nodes in $\mathcal{V}_t^{\text{sink}}$. We shall prove that the outgoing transmissions of every node in $\mathcal{V}^{\text{sink}}$, including those of V_{op} , must be independent of the message, i.e.,

$$I(M; X(\mathcal{Q}(V))) = 0 \quad \forall V \in \mathcal{V}^{\text{sink}}. \quad (40)$$

Initialization. Let t_0 be the first time instant t for which $\mathcal{V}_t^{\text{sink}}$ is non-empty. Then, we encounter two cases: either $t_0 = 0$, in which case the nodes in $\mathcal{V}_{t_0}^{\text{sink}}$ have *no* incoming edges, or $t_0 > 0$, and the nodes in $\mathcal{V}_{t_0}^{\text{sink}}$ *have* incoming edges. We shall first prove that in *both* cases, the outgoing transmissions of $\mathcal{V}_{t_0}^{\text{sink}}$ are independent of the message, i.e. $I(M; X(\mathcal{Q}(\mathcal{V}_{t_0}^{\text{sink}}))) = 0$.

(Case I) When $t_0 = 0$, $\mathcal{V}_0^{\text{sink}} \cap \mathcal{V}_{\text{ip}} = \emptyset$. This is because the cut separates \mathcal{V}_{ip} from V_{op} , with $\mathcal{V}_{\text{ip}} \subseteq \mathcal{V}^{\text{src}}$, so no nodes in $\mathcal{V}_0^{\text{sink}}$ can be input nodes. So, by the definition of (non-)input nodes (Definition 3c), we must have

$$I(M; X(\mathcal{Q}(\mathcal{V}_0^{\text{sink}}))) = I(M; f_{\mathcal{V}_0^{\text{sink}}}(W(\mathcal{V}_0^{\text{sink}}))) \quad (41)$$

$$\stackrel{(a)}{\leq} I(M; W(\mathcal{V}_0^{\text{sink}})) \quad (42)$$

$$\stackrel{(b)}{=} 0, \quad (43)$$

where step (a) uses the data processing inequality and step (b) makes use of the fact that $W(\mathcal{V}_0) \perp M$.

(Case II) When $t_0 > 0$, the definition of t_0 implies that all nodes at time $t_0 - 1$ are in \mathcal{V}^{src} , so all incoming edges of $\mathcal{V}_{t_0}^{\text{sink}}$ must lie in the cut set, i.e., $\mathcal{P}(\mathcal{V}_{t_0}^{\text{sink}}) \subseteq \mathcal{E}^{\text{cut}}$. Since the cut is a zero- M -information cut, we have that for all $E_{t_0-1} \in \mathcal{P}(\mathcal{V}_{t_0}^{\text{sink}})$,

$$I(M; X(E_{t_0-1}) | X(\mathcal{E}'_{t_0-1})) = 0 \quad \forall \mathcal{E}'_{t_0-1} \subseteq \mathcal{E}_{t_0-1}. \quad (44)$$

By the definition of M -information flow for a set of edges (Definition 5) and Proposition 1, we have

$$I(M; X(\mathcal{P}(\mathcal{V}_{t_0}^{\text{sink}})) \mid X(\mathcal{E}'_{t_0-1})) = 0 \quad \forall \mathcal{E}'_{t_0-1} \subseteq \mathcal{E}_{t_0-1}. \quad (45)$$

Once again, considering $\mathbb{Q}(\mathcal{V}_{t_0}^{\text{sink}})$, we have

$$I(M; X(\mathbb{Q}(\mathcal{V}_{t_0}^{\text{sink}}))) = I(M; f_{\mathcal{V}_{t_0}^{\text{sink}}}(X(\mathcal{P}(\mathcal{V}_{t_0}^{\text{sink}})), W(\mathcal{V}_{t_0}^{\text{sink}}))) \quad (46)$$

$$\stackrel{(a)}{\leq} I(M; X(\mathcal{P}(\mathcal{V}_{t_0}^{\text{sink}})), W(\mathcal{V}_{t_0}^{\text{sink}})) \quad (47)$$

$$\stackrel{(b)}{=} I(M; X(\mathcal{P}(\mathcal{V}_{t_0}^{\text{sink}}))) + I(M; W(\mathcal{V}_{t_0}^{\text{sink}}) \mid X(\mathcal{P}(\mathcal{V}_{t_0}^{\text{sink}}))) \quad (48)$$

$$\stackrel{(c)}{=} 0, \quad (49)$$

where (a) and (b) follow from the Data Processing Inequality and the chain rule of mutual information respectively. In step (c), the first expression in the sum goes to zero by taking $\mathcal{E}_{t_0-1} = \emptyset$ in (45) and the second expression is zero since $W(\mathcal{V}_{t_0}^{\text{sink}}) \perp \{M, X(\mathcal{E}_{t_0-1})\}$, and $\mathcal{P}(\mathcal{V}_{t_0}^{\text{sink}}) \subseteq \mathcal{E}_{t_0-1}$ (refer Definition 3b). So, from equations (43) and (49), we have that for all values of t_0 ,

$$I(M; X(\mathbb{Q}(\mathcal{V}_{t_0}^{\text{sink}}))) = 0. \quad (50)$$

Continuation. Now, suppose that for some $t > t_0$, we have $I(M; X(\mathbb{Q}(\mathcal{V}_{t-1}^{\text{sink}}))) = 0$. We shall prove that this implies $I(M; X(\mathbb{Q}(\mathcal{V}_t^{\text{sink}}))) = 0$. First, observe that

$$\mathcal{P}(\mathcal{V}_t^{\text{sink}}) = (\mathcal{P}(\mathcal{V}_t^{\text{sink}}) \cap \mathcal{E}^{\text{cut}}) \cup (\mathcal{P}(\mathcal{V}_t^{\text{sink}}) \cap \mathcal{E}^{\text{sink}}) \quad (51)$$

For convenience, let $\mathcal{P}_{t-1}^{\text{cut}} := \mathcal{P}(\mathcal{V}_t^{\text{sink}}) \cap \mathcal{E}^{\text{cut}}$ and $\mathcal{P}_{t-1}^{\text{sink}} := \mathcal{P}(\mathcal{V}_t^{\text{sink}}) \cap \mathcal{E}^{\text{sink}}$. We have used the subscript $t-1$ here to remind the reader that $\mathcal{P}(\mathcal{V}_t^{\text{sink}})$, which are the *incoming* edges of $\mathcal{V}_t^{\text{sink}}$, are a subset of \mathcal{E}_{t-1} . Then, we have

$$\mathcal{P}(\mathcal{V}_t^{\text{sink}}) = \mathcal{P}_{t-1}^{\text{cut}} \cup \mathcal{P}_{t-1}^{\text{sink}}. \quad (52)$$

Since the cut is a zero- M -information cut, we have that for every $E_{t-1} \in \mathcal{P}_{t-1}^{\text{cut}}$,

$$I(M; X(E_{t-1}) \mid X(\mathcal{E}'_{t-1})) = 0 \quad \forall \mathcal{E}'_{t-1} \subseteq \mathcal{E}_{t-1}. \quad (53)$$

Therefore, by Definition 5 and Proposition 1,

$$I(M; X(\mathcal{P}_{t-1}^{\text{cut}}) \mid X(\mathcal{E}'_{t-1})) = 0 \quad \forall \mathcal{E}'_{t-1} \subseteq \mathcal{E}_{t-1}. \quad (54)$$

Secondly, $\mathcal{P}_{t-1}^{\text{sink}} = \mathbb{Q}(\mathcal{V}_{t-1}^{\text{sink}}) \cap \mathcal{E}^{\text{sink}}$. This is depicted in Figure 5, and explained in the caption. So,

$$I(M; X(\mathcal{P}_{t-1}^{\text{sink}})) = I(M; X(\mathbb{Q}(\mathcal{V}_{t-1}^{\text{sink}}) \cap \mathcal{E}^{\text{sink}})) \quad (55)$$

$$\stackrel{(a)}{\leq} I(M; X(\mathbb{Q}(\mathcal{V}_{t-1}^{\text{sink}}))) \stackrel{(b)}{=} 0 \quad (56)$$

where (a) follows from the fact that considering more random variables can only increase mutual information, and (b) follows from the induction assumption. Finally, consider how $X(\mathbb{Q}(\mathcal{V}_t^{\text{sink}}))$ depends on M :

$$I(M; X(\mathbb{Q}(\mathcal{V}_t^{\text{sink}}))) = I(M; f_{\mathcal{V}_t^{\text{sink}}}(X(\mathcal{P}_{t-1}^{\text{sink}} \cup \mathcal{P}_{t-1}^{\text{cut}}), W(\mathcal{V}_t^{\text{sink}}))) \quad (57)$$

$$\stackrel{(a)}{\leq} I(M; X(\mathcal{P}_{t-1}^{\text{sink}}), X(\mathcal{P}_{t-1}^{\text{cut}}), W(\mathcal{V}_t^{\text{sink}})) \quad (58)$$

$$\stackrel{(b)}{=} I(M; X(\mathcal{P}_{t-1}^{\text{sink}})) + I(M; X(\mathcal{P}_{t-1}^{\text{cut}}) \mid X(\mathcal{P}_{t-1}^{\text{sink}})) \quad (59)$$

$$+ I(M; W(\mathcal{V}_t^{\text{sink}}) \mid X(\mathcal{P}_{t-1}^{\text{sink}}), X(\mathcal{P}_{t-1}^{\text{cut}})) \quad (60)$$

$$\stackrel{(c)}{=} 0,$$

where once again, (a) and (b) follow from the data processing inequality and the chain rule respectively. In step (c), the first and second terms go to zero by equations (56) and (54) respectively, while the third term is zero since $W(\mathcal{V}_t^{\text{sink}}) \perp \{M, X(\mathcal{E}_{t-1})\}$ and $\mathcal{P}_{t-1}^{\text{sink}} \cup \mathcal{P}_{t-1}^{\text{cut}} \subseteq \mathcal{E}_{t-1}$.

The proof follows from induction on t , so

$$I(M; X(\mathbb{Q}(\mathcal{V}_t^{\text{sink}}))) = 0 \quad \forall t \geq t_0, \quad (61)$$

which in turn implies that

$$I(M; X(\mathbb{Q}(V))) = 0 \quad \forall V \in \mathcal{V}^{\text{sink}}. \quad (62)$$

If there exists an output node whose transmissions depend on M , then there can exist no cut consisting of edges with zero M -information flow, and hence by Lemma 8, there must be a path consisting of edges that carry M -information flow between the input nodes and the output node in question. \square

4.4 The Separability Property

Finally, we state a property that may be of interest to obtain a deeper understanding of the nature of M -information flow, as given by Definitions 4 and 5.

Proposition 9 (Separability). *Let \mathcal{C} be a computational system. Then, at any given point in time t , there exist two sets $\mathcal{R}_t, \mathcal{S}_t \subseteq \mathcal{E}_t$, such that all of the following conditions hold:*

1. $\mathcal{R}_t \cup \mathcal{S}_t = \mathcal{E}_t$
2. $\mathcal{R}_t \cap \mathcal{S}_t = \emptyset$
3. Either $\mathcal{R}_t = \emptyset$, or for every $R_t \in \mathcal{R}_t$ there exists a subset $\mathcal{R}'_t \subseteq \mathcal{R}_t \setminus \{R_t\}$ such that

$$I(M; X(R_t) \mid X(\mathcal{R}'_t)) > 0. \quad (63)$$

4. Either $\mathcal{S}_t = \emptyset$, or for every $\mathcal{E}'_t \subseteq \mathcal{E}_t$,

$$I(M; X(\mathcal{S}_t) \mid X(\mathcal{E}'_t)) = 0. \quad (64)$$

A proof of this proposition can be found in Appendix B.

Proposition 9 shows that at any given point in time t , it is possible to partition \mathcal{E}_t into two sets: \mathcal{R}_t , consisting only of edges that have M -information flow, and \mathcal{S}_t , comprising edges that have no M -information flow. Furthermore, when considering the M -information flow of edges in \mathcal{R}_t , it suffices to condition on the transmissions of edges within \mathcal{R}_t to ascertain the presence of M -information flow. Conditioning upon the transmissions of edges in \mathcal{S}_t will not change the mutual information between the message and the transmissions of edges in \mathcal{R}_t .

5 Inferring Information Flow

Having discussed the definition and the properties of M -information flow, we now consider how these flows of information might be inferred in a real computational system. We first discuss an observation model that describes which random variables are observed and how they are sampled. Under this model, we show how existing techniques from the literature can be used to identify which edges carry M -information flow. As in previous sections, we restrict our attention to detecting *whether or not* a given edge has M -information

flow, relegating quantification of these flows to future work. Quantification is briefly discussed in the form of an example in Section 6.3, and again in Section 7.5.

We then describe an algorithm that recovers all M -information paths between the input nodes and a given output node, by leveraging the knowledge of which edges have M -information flow. We also explain how one might attain a fine-grained characterization of the structure of information flow, by introducing the concept of “derived information”. This is useful for understanding which transmissions are “derived” from others, allowing one to find transmissions that are redundant and discover the presence of hidden nodes. Finally, we explain how flows of information about multiple messages can be inferred in our framework.

5.1 The Observation Model

Before we can describe how information flow and information paths can be identified, we must provide a statistical description of the random variables that are observed. Let \mathcal{C} be a computational system under observation. We then make the following assumptions:

1. Transmissions on all edges, including self-edges, are observed. The random variables that are intrinsically generated at each node are *not* observed, unless they are also transmitted on an edge (which could be a self-edge).
2. Several trials²¹ are observed, each of which corresponds to an independent realization of all random variables in the model²². Every trial uses a realization of M which is independently drawn from a distribution determined by the experimentalist²³. For every node $V \in \mathcal{V}$, the intrinsically generated random variable $W(V)$ is also assumed to be independently and identically distributed across trials.
3. Observations are made noiselessly, in that the realization of each transmission in every trial is observed as-is, without being further corrupted by random noise of any kind. The implications of noisy measurements will be the subject of future work.

Under these conditions, we discuss statistical tests for information flow that are consistent in the asymptotic limit of infinite trials. It should be noted that these assumptions may be valid to varying degrees in different contexts. This is discussed further in Section 7.1.

5.2 Detecting Information Flow

Given a sample of all random variables described in the observation model, our next task is to identify which edges have M -information flow. In other words, we need to describe how the conditions given by Definition 4 can be rigorously tested, and how we might assert with some confidence that a certain set of edges has information flow at each point in time.

According to Definition 4, in order to check whether a particular edge E_t carries M -information flow at time t , we need to test whether at least one of several conditional mutual information quantities is strictly positive.

²¹The word “trial” is borrowed from the neuroscience literature, wherein a neuroscientist will often conduct multiple trials in a single experiment. In each trial, a human participant or an animal under study is presented with one of a set of carefully chosen stimuli (corresponding to a realization of the message M in our setting), and neural activity is recorded using some modality. Scientific inferences are then drawn by making use of the activity from all trials.

²²In reality, trials are not independent in neuroscientific experiments. Indeed, neurons are known to “adapt” their responses from trial to trial, often showing suppressed activity when presented the same stimulus multiple times. This, in part, is considered to be evidence of *learning* in neural circuitry. However, for simplicity, we restrict our attention here to computational systems that do not learn or show trial-to-trial adaptation.

²³A more detailed discussion of this distribution can be found in Section 5.6.

The standard statistical approach for solving this problem is to frame it as a set of “hypothesis tests”, which in this case is a set of “conditional independence tests”. In general, a hypothesis test formalizes the problem of making an informed decision about the value of some functional of a joint distribution, when observing a sample of data from it. A good conditional independence testing procedure will seek to maximize “statistical power”, i.e. the probability of *correctly* identifying the presence of conditional dependence, while keeping the probability of an incorrect identification fixed below some “level” α that is picked beforehand. One intuitive way to do this might be to construct an estimator for the appropriate conditional mutual information, and “reject” the “null” hypothesis of conditional independence if the conditional mutual information was sufficiently larger than some threshold, $\epsilon > 0$. This threshold would have to be chosen so that, on average, the probability of falsely rejecting the null hypothesis is at most α . However, there are usually better ways of performing this test, i.e., it is often possible to attain higher power at the same level *without* actually estimating the conditional mutual information.

While it would be impossible to provide a comprehensive list of papers that have researched the problem of conditional independence testing, it has received (and continues to receive) much attention in the statistics, causality, and information theory communities [57–62]. In its most general form, conditional independence testing is considered to be a hard problem for continuous random variables [63]. However, if we ignore issues associated with the practical difficulty of estimation (discussed later in Section 7.2), these works provide consistent tests under reasonable assumptions on the joint distribution of the variables involved [59–61].

Although we mentioned that there are better ways to test for conditional dependence than to estimate the conditional mutual information, there may be instances when one might want to estimate the conditional mutual information anyway. For instance, in an example that will appear shortly in Section 6.3, we rely on an *estimate* of the conditional mutual information to *quantify* the amount of M -information flowing on a given edge. While our paper has only defined M -information flow in terms of *whether or not* it is present at an edge E_t , it is also extremely useful to know *how much* M -information flow there is. We defer further discussion of this topic until Sections 6.3 and 7.5. For now, we note that several papers have considered how to *estimate* mutual information and conditional mutual information, both of which might be essential for an understanding of *quantification* of M -information flow [64–67].

For completeness, we now present a description of how we expect information flow will be detected in practice. We assume that we have samples of observations from every edge of the computational system, at every point in time. If not, appropriate assumptions may need to be made, as discussed later in Section 7.1. At every instant of time t , consider the set of all edges \mathcal{E}_t present in the network. For every edge $E_t \in \mathcal{E}_t$, use the following process to determine whether it has M -information flow:

1. First test whether the mutual information between its transmission and the message is greater than zero, i.e., $I(M; X(E_t)) > 0$. If so, declare that E_t has M -information flow.
2. If not, test for conditional dependence between its transmission and the message, given each of the other edges E'_t , i.e., $I(M; X(E_t) | X(E'_t)) > 0, \forall E'_t \in \mathcal{E}_t \setminus \{E_t\}$. If any of these tests rejects the null hypothesis, declare that E_t has M -information flow.
3. If not, test for conditional dependence between $X(E_t)$ and M , given subsets of other edges, while sequentially taking edges taken pairwise, then in threes, etc. If any of these tests rejects the null, declare that E_t has M -information flow.
4. If none of the above tests rejects the null hypothesis, declare that E_t carries no M -information flow.

Note that we have not discussed the level, α , at which we should reject the null in each of the above tests. In general, since we are performing multiple hypothesis tests simultaneously, some manner of “*correction*” is required to ensure that we do not find, what is effectively, a spurious correlation. This is discussed at length in Section 7.2.

Algorithm 1 Information Path Algorithm: Finds all paths from \mathcal{V}_{ip} to V_{op}

```

1: Initialize an empty graph  $\mathcal{H}$  ▷  $\mathcal{H}$  will store valid paths from  $\mathcal{V}_{ip}$  to  $V_{op}$ 
▷  $\mathcal{H}$  currently contains no nodes or edges
2: FINDINFOPATHS( $\mathcal{C}$ ,  $V_{op}$ ,  $\mathcal{H}$ ) ▷ Call a function (defined below) to populate  $\mathcal{H}$ 
3: if  $V_{op}$  is marked “invalid” then
4:   raise Error ▷ No path from  $\mathcal{V}_{ip}$  to  $V_{op}$  was found
5: end if

6: function FINDINFOPATHS( $\mathcal{C}$ ,  $V_t$ ,  $\mathcal{H}$ )
7:   if  $\mathcal{P}(V_t)$  is empty then ▷  $V_t$  has no inputs  $\Rightarrow t = 0$ 
8:     if  $V_t \in \mathcal{V}_{ip}$  then
9:       Mark  $V_t$  “valid”
10:      Add  $V_t$  to  $\mathcal{H}$ 
11:     else ▷ We somehow reached a non-input node at  $t = 0$ 
12:       raise Error
13:     end if
14:   else ▷  $V_t$  has inputs
15:     for all  $(U_{t-1}, V_t) \in \mathcal{P}(V_t)$  do
16:       if  $(U_{t-1}, V_t)$  has  $M$ -information flow then
17:         if  $U_{t-1}$  is unmarked then
18:           FINDINFOPATHS( $\mathcal{C}$ ,  $U_{t-1}$ ,  $\mathcal{H}$ ) ▷ This will mark  $U_{t-1}$ 
19:         end if
20:         if  $U_{t-1}$  is marked “valid” then
21:           Mark  $V_t$  “valid”
22:           Add  $V_t$  and  $(U_{t-1}, V_t)$  to  $\mathcal{H}$ 
23:         end if
24:       end if
25:     end for
26:     if  $V_t$  is still unmarked then ▷ No input of  $V_t$  was “valid”
27:       Mark  $V_t$  “invalid”
28:     end if
29:   end if
30: end function

```

5.3 Discovering Information Paths

Next, we discuss an algorithm that discovers all M -information paths leading from the input nodes to a given output node, V_{op} , in any computational system. As discussed in Section 4, whenever the transmissions $\mathbb{Q}(V_{op})$ of the output node depend on the message, Theorem 7 guarantees that at least one M -information path exists.

Algorithm 1, which we propose for recovering all M -information paths, is an adaptation of the well-known Depth-First Search²⁴ method [56, Sec. 22.3]. It takes as its input a computational system \mathcal{C} in which all edges having M -information flow have been identified, the output node V_{op} , and an empty graph \mathcal{H} that is completely devoid of nodes and edges. The algorithm returns the set of all M -information paths in the form of a directed subgraph \mathcal{H} of the time-unrolled graph \mathcal{G} . Starting from \mathcal{V}_{ip} , following *any* path in \mathcal{H} will lead

²⁴It is also possible to discover all M -information paths using an adaptation of Breadth-First Search [56, Sec. 22.2], but doing so would require some mechanism to prune M -information paths that do not lead to the input nodes \mathcal{V}_{ip} . So we prefer to use Depth-First Search for simplicity of exposition.

one to V_{op} , provided at least one M -information path exists.

The algorithm works by recursively visiting nodes, starting from the output node V_{op} . It traverses only edges that carry M -information flow, and uses a marking scheme to avoid revisiting nodes. The same marking scheme is also used to designate nodes to which there are M -information paths from \mathcal{V}_{ip} . As the algorithm passes through each node, it marks the node “valid” whenever an M -information path exists between \mathcal{V}_{ip} and that node. If no such path exists, then the node is marked “invalid”. The objective of the algorithm, therefore, reduces to one of finding a path of “valid” nodes from \mathcal{V}_{ip} to V_{op} . The algorithm’s recursive function can be expressed as follows: *A node $V_t \in \mathcal{V}$ is “valid” if and only if there exists a node $U_{t-1} \in \mathcal{V}$ such that U_{t-1} is valid, and the edge (U_{t-1}, V_t) has M -information flow.* This is a recursive expression since checking the validity of a node at time t involves finding valid nodes at time $t - 1$. The only nodes that are considered valid by default are the input nodes \mathcal{V}_{ip} .

The algorithm sequentially checks the validity of nodes $V_t \in \mathcal{V}$, starting from the output node V_{op} . The function `FINDINFOPATHS`, when called on any given node V_t , checks the validity of V_t . This involves checking each of the incoming edges of V_t for M -information flow. If U_{t-1} is a node from which M -information flows to V_t , then the algorithm immediately checks the validity of U_{t-1} by calling the function `FINDINFOPATHS` again. Eventually, if in this recursive process, we arrive at an input node in \mathcal{V}_{ip} , then that node is marked “valid”, and added to the output subgraph \mathcal{H} . Once every node U_{t-1} from which M -information flows to V_t has been marked “valid” or “invalid”, the validity of V_t can be ascertained. For every “valid” node U_{t-1} from which M -information flows to V_t , the edge (U_{t-1}, V_t) and the node V_t are added to the output subgraph \mathcal{H} , and V_t is marked “valid”. If there are no such nodes leading to V_t , then V_t is marked “invalid” and does not fall on an M -information path.

This recursive logic yields the set of all M -information paths leading from the input nodes to V_{op} . The two lines at which errors are returned correspond to scenarios that should not occur if the conditions of Theorem 7 hold. In line 12, we visit a non-input node at time $t = 0$. But such a node should never have been reached in the recursion, since we only followed edges that have M -information flow. Its presence, therefore, would contradict the computational system model. In line 4, V_{op} is marked “invalid”, implying that there is no path leading to it from the input nodes. Once again, this can only occur if the computational system model is violated, or if the conditions of Theorem 7 do not hold.

On Computational Complexity

The complexity of this algorithm is exactly that of Depth-first Search, $\mathcal{O}(|\mathcal{V}| + |\mathcal{E}|)$ [56, Sec. 22.3]. To be precise, we consider the computational system to extend until the time of the output node, i.e., we take $T = t_{\text{op}}$. So the complexity of the algorithm is $\mathcal{O}(|\mathcal{V}^*|t_{\text{op}} + |\mathcal{E}^*|t_{\text{op}})$. This is easily verified: if we assume that all edges in the system have M -information flow, then all edges and nodes must be traversed by the search. At each node, we must execute lines 7 through 14, and 26 through 28, which take a constant amount of time. Since we have $|\mathcal{V}^*|$ nodes over t_{op} time points, this adds up to $\mathcal{O}(|\mathcal{V}^*|t_{\text{op}})$ steps. We also need to execute the loop in lines 15 through 24, which counts the number of incoming edges at every node. For all nodes combined, this adds up to $\mathcal{O}(|\mathcal{E}^*|t_{\text{op}})$ steps.

If the graph is fully connected as described in Section 2, then $|\mathcal{V}^*| = N$ and $|\mathcal{E}^*| = N^2$, so the effective complexity is just $\mathcal{O}(N^2t_{\text{op}})$. However, if we *know* that the underlying graph is sparse (e.g., because of anatomical priors in neuroscience), then we may have $|\mathcal{E}^*| = \mathcal{O}(N \log N)$, or even $|\mathcal{E}^*| = \mathcal{O}(N)$, bringing down the complexity of the search. It should be noted that in either case, the complexity of identifying *which* edges have M -information flow is potentially exponential in N , as discussed later in Section 7.2. This is much larger than the complexity of tracing out information paths, so *finding edges* with M -information flow is, in fact, the “hard part” of the problem.

5.4 Derived Information and Redundancy

The framework we develop for information flow allows one to obtain a more fine-grained understanding of information structure in a computational system, especially when compared with classical tools such as correlation and phase synchrony [68, 69]. This allows the experimentalist to better investigate the nature of the computation being performed. A concept that we believe will be extremely useful in this regard is one we call “derived information”, which is defined below.

Definition 12 (Derived M -Information). *In a computational system \mathcal{C} , a transmission $X(Q_t)$ is said to be derived M -information of a different transmission $X(P_{t'})$ if $M - X(P_{t'}) - X(Q_t)$ forms a Markov chain. That is, the following condition must hold:*

$$I(M; X(Q_t) \mid X(P_{t'})) = 0, \quad (65)$$

implying that

$$H(M \mid X(P_{t'})) = H(M \mid X(P_{t'}), X(Q_t)). \quad (66)$$

So, $X(Q_t)$ adds no new information about M , when given $X(P_{t'})$. The same definition extends to transmissions on sets of edges. Note that, as far as the definition is concerned, t and t' may be any two arbitrary points in time. However, we will typically consider cases when $t \geq t'$.

One potential use-case scenario for derived information arises in the context of redundant flows. Consider the computational system presented in Figure 4, originally described under Counterexample 3. We see two edges sending the same transmission to the node B_2 . This is an example of what we call “redundant transmissions”. In general, since we only consider information about M to be relevant, the exact transmissions communicated over two edges at a given point in time may be different. But if they convey the *same information about M* to a given node, then we view them as essentially redundant. Definition 4, when applied to this system, will detect both these edges as having M -information flow, since given $X((C_1, B_2))$, their transmissions depend on M . In the notation of the Separability property mentioned earlier (Proposition 9), both edges (A_1, B_2) as well as (D_1, B_2) will belong in the set \mathcal{R}_1 .

Derived information provides a general methodology to understand when transmissions on certain edges may be redundant. Naturally, if the transmissions on two edges Q_t and $P_{t'}$ are redundant, then they must be derived M -information of one another. This amounts to checking two more conditional independence relationships, for which consistent tests exist in the limit of infinite trials, as discussed in Section 5.2.

In the following section, we shall see another application of derived information; when applied to specific sets, it can in some cases be used to detect the presence of hidden (unobserved) nodes. Later, in Section 6.3, we discuss an example where the notion of derived information helps us make a new kind of inference about the fine structure of information flow, one that would not be possible using tools such as Granger Causality and Directed Information.

5.5 Hidden Nodes

In Section 5.3, we showed how the Information Path Algorithm may fail to discover M -information paths if one of the assumptions of the computational system model or the observation model breaks in some way. Here, we discuss one specific situation in which the observation model may break, i.e., when not all nodes are observed. We call these unobserved nodes “hidden nodes”, and assume that we do not see transmissions on incoming or outgoing edges of these nodes.

Definition 13 (Hidden nodes). *Consider a computational system $\mathcal{C} = (\mathcal{G}, X, W, f)$ defined on the time-unrolled graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ as before. Suppose that only a subset of nodes in this graph are observed.*

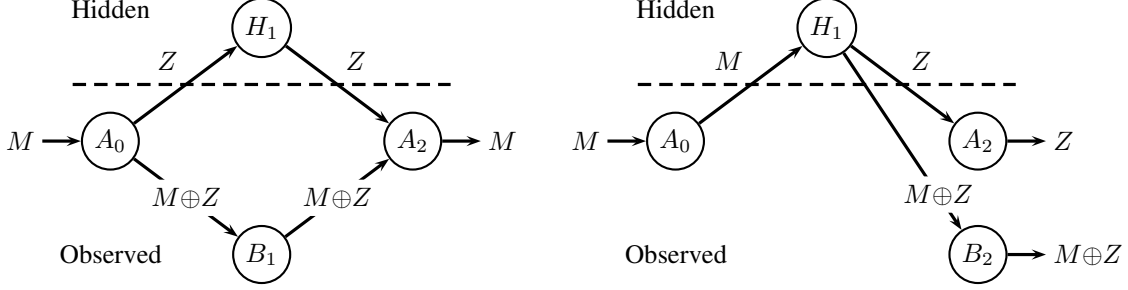


Figure 6: Two simple examples showing how hidden nodes may prevent one from being able to discover M -information paths in a computational system. In both cases shown here, H_1 is a hidden node, and we do not observe its incoming or outgoing transmissions. On the left is an example where a transmission that we might need to condition upon to discover M -information flow passes through the hidden node, and therefore cannot be seen. On the right, the hidden node itself generates the source of randomness Z .

Specifically, if \mathcal{V}^* was the original set of nodes in \mathcal{G}^* , prior to time-unrolling, then we observe only the nodes $\tilde{\mathcal{V}}^* = \mathcal{V}^* \setminus \mathcal{H}^*$, where $\mathcal{H}^* = \{H^{(0)}, H^{(1)}, \dots, H^{(K-1)}\}$ is a set of unobserved nodes called hidden nodes.

To describe the observed component of the computational system, we define $\tilde{\mathcal{G}}^* = \tilde{\mathcal{V}}^* \times \tilde{\mathcal{V}}^*$, $\tilde{\mathcal{V}} = \{V_t : V \in \tilde{\mathcal{V}}^*, t \in \mathcal{T}\}$ and $\tilde{\mathcal{E}} = \{(A_t, B_{t+1}) : (A, B) \in \tilde{\mathcal{G}}^*, t \in \mathcal{T}\}$. Also let $\mathcal{H} = \{H_t : H \in \mathcal{H}^*, t \in \mathcal{T}\}$. Finally, we set up the observed component of the computational system as before: $\tilde{\mathcal{G}} = (\tilde{\mathcal{G}}, X, W, f)$. Thus, we only observe the transmissions on edges in $\tilde{\mathcal{E}}$. As usual, we denote the set of all hidden nodes at time t by \mathcal{H}_t , and the set of all observed nodes at time t by $\tilde{\mathcal{E}}_t$.

The presence of hidden nodes of this nature implies that much of the theory we have developed will not apply. Lemma 4 no longer truly holds, in that information about M may persist in the system by passing through the hidden node, even if *no* observed edge has M -information flow. So, naturally, Property 1 also fails to hold. Hence, we are not guaranteed to be able to identify all edges with M -information flow, and discover all M -information paths as before. For example, refer to the cases shown in Figure 6, where we no can longer find M -information paths because of the presence of a hidden node.

Fortunately, at least in some cases, the concept of derived information (Definition 12) provides a simple way to tell whether or not a hidden node exists. Specifically, if at some time t , a hidden node transmits information about M which is unavailable within the system at that time, and which is utilized by some node at the next time instant, then the set of all observed transmissions $X(\tilde{\mathcal{E}}_t)$ will *not* be derived M -information of the set of all transmissions at time $t-1$. In other words, the Global Markovity condition (Corollary 5) on the observed graph, $M-X(\tilde{\mathcal{E}}_{t-1})-X(\tilde{\mathcal{E}}_t)$, will break. Unfortunately, the notion of “utilization” is difficult to express mathematically, without resorting to the use of ideas from causality that are based on intervention. The result we prove, therefore, is a simpler sufficiency argument, which guarantees the presence of a hidden node if the aforementioned Markov condition is observed to break. This result is proved in Proposition 11, but first, we define some adjectives.

Definition 14 (M -relevant hidden node). *A hidden node H_t is said to be M -relevant if $\mathbb{Q}(H_t)$ carries M -information flow in \mathcal{G} . Similarly, a subset of hidden nodes $\mathcal{H}'_t \subseteq \mathcal{H}_t$ is said to be M -relevant if $\mathbb{Q}(\mathcal{H}'_t)$ carries M -information flow in \mathcal{G} .*

Definition 15 (M -derived hidden node). *A hidden node H_t is said to be M -derived if the Markov chain $M-X(\tilde{\mathcal{E}}_t)-X(\mathbb{Q}(H_t))$ holds. Similarly, a subset of hidden nodes $\mathcal{H}'_t \subseteq \mathcal{H}_t$ is said to be M -derived if the Markov chain $M-X(\tilde{\mathcal{E}}_t)-X(\mathbb{Q}(\mathcal{H}'_t))$ holds.*

Lemma 10. *If a subset of hidden nodes is not M -derived, then it is M -relevant.²⁵*

²⁵If this lemma appears to be somewhat strong, it is only because of the nomenclature “ M -derived”. For our purposes, a

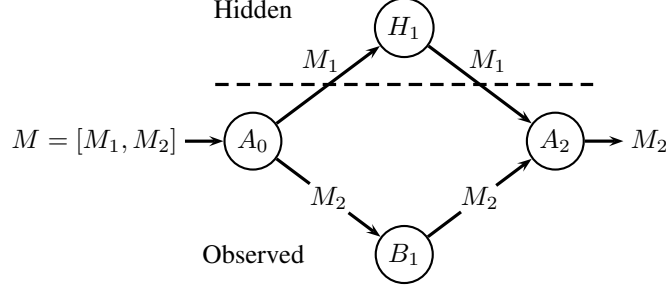
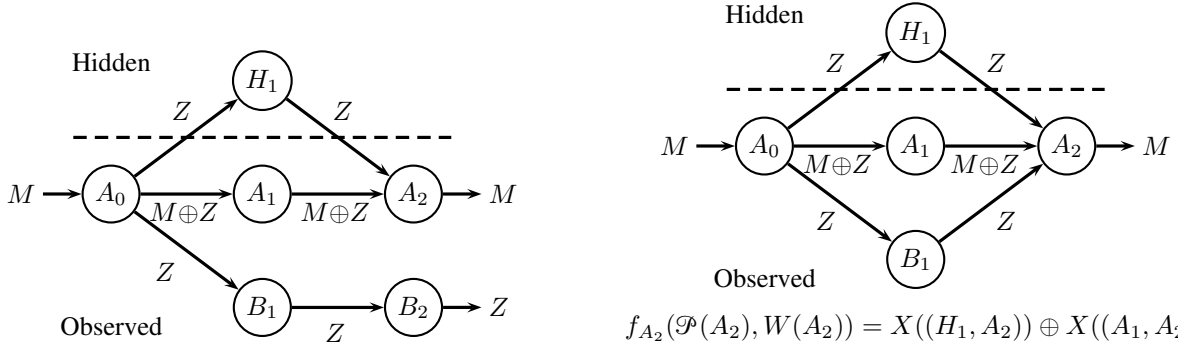


Figure 7: A computational system serving as a counterexample to the converse of Proposition 11. Here, the hidden node H_1 is M -relevant because its outgoing transmission, M_1 , is not present in any of the observed transmissions at time $t = 1$. However, since A_2 chooses to ignore M_1 at its output, the Markov chain $M - X(\mathcal{E}_1) - X(\mathcal{E}_2)$ boils down to $M - M_2 - M_2$, which obviously holds. Thus, at least based on our current definitions, there may be M -relevant hidden nodes in the system even if Global Markovity continues to hold.



(a) Here, although Global Markovity holds, one could argue that testing for Local Markovity at each node (or at various subsets of nodes) could help uncover the presence of a hidden node.

$$f_{A_2}(\mathcal{P}(A_2), W(A_2)) = X((H_1, A_2)) \oplus X((A_1, A_2))$$

(b) In this case, the hidden node breaks neither Global nor Local Markovity. However, the *function* computed by A_2 makes use of only the hidden node's transmission. As a result, the hidden node has a causal effect on the output of the system, since destroying the outgoing edge of the hidden node would change the output. Such a hidden node is likely undetectable using only observational methods.

Figure 8: Examples of computational systems with an M -derived hidden node. In both of these systems, the hidden node's transmission at time $t = 1$ has an affect on the output at A_2 . However, Global Markovity continues to hold from $t = 1$ to $t = 2$, because the observed transmissions, $M \oplus Z$ and Z , contain all information necessary to explain the output, M .

Proposition 11. *In a computational system \mathcal{C} with hidden nodes, if Global Markovity on the observed graph, $\tilde{\mathcal{G}}$, fails to hold from time t to $t + 1$, i.e. if $I(M; X(\tilde{\mathcal{E}}_{t+1}) | X(\tilde{\mathcal{E}}_t)) > 0$, then the hidden nodes \mathcal{H}_t at time t are not M -derived.*

Proofs of Lemma 10 and Proposition 11 are very straightforward, and are provided in Appendix D. As a direct consequence of these two results, if Global Markovity fails to hold on the observed nodes from time t to $t + 1$, then \mathcal{H}_t is M -relevant. By Proposition 1, this simply means that there exists at least one M -relevant hidden node at time t .

Although Proposition 11 appears to provide a straightforward mechanism to test whether or not hidden nodes exist, it does not always work. If a hidden node's transmissions have no M -information flow, then the node will not be detected. But in this case, it could be argued that such a hidden node does not change whether information paths can be identified, and so can be subsumed by one or more of the intrinsic random

hidden node whose transmissions are independent of the message is also M -derived, since it satisfies the aforementioned Markov condition.

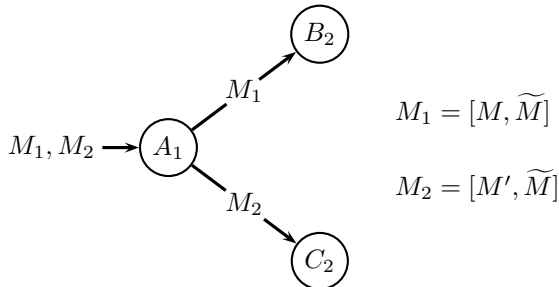


Figure 9: A simple example demonstrating the importance of having independent messages (or sub-messages) when exploring the flows of multiple messages in a computational system. As M_1 and M_2 both redundantly contain information about \tilde{M} , both edges shown here have M_1 - as well as M_2 -information flow. Thus, we are unable to detect the fact that M_1 and M_2 take different paths in the system, because of our choice of stimuli.

variables $W(\cdot)$. Such a hidden node is, therefore, classified by Definition 14 as *not* M -relevant.

However, to make matters worse, the converse of Proposition 11 does not hold. In particular, there may exist an M -relevant hidden node at time t , whose transmission is *ignored* by the node that received it, so that the Markov chain $M-X(\tilde{\mathcal{E}}_t)-X(\tilde{\mathcal{E}}_{t+1})$ continues to hold (see Figure 7). Such a hidden node may *still* be considered largely innocuous.

The most serious case of a hidden node going undetected is one that contains an M -derived hidden node, whose transmission *is used* by the receiving node while performing its computation; however the hidden node’s transmission is “masked” by a redundant transmission from an observed node (see Figure 8). In this case, Global Markovity on $\tilde{\mathcal{G}}$ will *not* break, yet the hidden node’s transmission may be instrumental in producing a certain output distribution. In some instances, such hidden nodes can be detected by checking for Local Markovity (Proposition 6; see Figure 8a). However, there are still cases where if we were somehow able to intervene and delete the transmission of the hidden node, then the computational system’s output may not remain the same, despite the existence of a redundant transmission from an observed node (see Figure 8b). Indeed, the presence of redundancy in such a scenario does not guarantee that the computational system will actually leverage it.

5.6 On Multiple Messages and the Distribution of the Message

Just as we can infer information flow and information paths for a single message, we can examine the flows of multiple messages in the same computational system. Consider a case where we wish to understand the information flows of two messages, M_1 and M_2 . An neuroscientific example of this might be information flow about two independent components of a visual stimulus, e.g., shape and color (such as in [2]). If $M_1 \perp M_2$, then we could separately identify edges and paths that have M_1 -information flow and M_2 -information flow, by applying the theory and algorithm as-is for each message individually.

However, if the two messages are *dependent* on one another, one could end up confounding their information flows, based on how they depend on each other, and how the computational system’s transmissions carry their joint information. As a simple example, consider the system shown in Figure 9, where $M_1 = [M, \tilde{M}]$ and $M_2 = [M', \tilde{M}]$, with $M, M', \tilde{M} \sim \text{i.i.d. Ber}(1/2)$. Clearly, M_1 and M_2 both share some redundant information in \tilde{M} , and $I(M_1; M_2) = 1$ bit. Thus, we will see M_1 -information flow as well as M_2 -information flow on both edges, since the transmission of each edge E satisfies $I(M_i; X(E)) > 0$ for $i \in \{1, 2\}$.

Consider what this means for the aforementioned example of shape and color of a visual stimulus. If a neuroscientist expects that the information paths corresponding to shape and color in the brain are different from each other, what is the best way to design stimuli so as to bring out this difference? Suppose they

decided to present a total of four different stimuli, $M \in \{0, 1, 2, 3\}$, with two different shapes and two different colors. Let M_1 be the first bit of the binary representation of M , denoting shape, and M_2 be the second bit, denoting color. Now if the neuroscientist chose to present stimuli with a uniform distribution over M , i.e., if each shape-color combination was shown for one-quarter of all trials, then M_1 and M_2 would be independent of each other, and their individual flows could be tracked separately. However, if the neuroscientist chose to present the four possible stimuli with probabilities $\{1/2, 1/4, 1/8, 1/8\}$ respectively, then M_1 and M_2 are no longer independent of each other, and it may become hard to separate their individual flows as in the example in Figure 9.

These examples suggest that, when trying to understand the flows of different messages in a computational system, it helps if they are independent of one another. So from the perspective of experiment design in a neuroscientific context, it is often more sensible to design stimuli so that the two messages of interest are independent of one another. Even when considering a single message that takes one of several values, it becomes important to appropriately choose a distribution over these values to ensure that any sub-messages that are of interest remain independent of one another. This would allow the experimentalist to better understand how “independent dimensions” of the stimulus are processed in the brain.

However, there are also situations where the experimental paradigm necessitates a statistical distribution of stimuli that makes two sub-messages of interest dependent on one another. For instance, the Posner experimental paradigm for attention [70] only works when the proportion of “valid” trials (a certain *type* of trial specific to this paradigm) is roughly 70%. Similarly, during data preprocessing, it is common to discard trials that are excessively noisy, based on some predetermined metric: this process could skew the distribution of the message, even if the original distribution was uniform. If it is still of interest to understand the individual flows of sub-messages in this case, then a possible solution might then be to sub-select experimental trials in such a way as to keep the two sub-messages independent of one another.

6 Canonical Computational Examples

In this section, we provide a few canonical examples for computational systems from various contexts. In each case, we discuss what the message M is, and identify which edges carry M -information flow. We also explain how the path recovered by the information path algorithm might be the intuitive choice in each example.

6.1 The Butterfly Network from Network Coding

For our first example, we cover the butterfly network from network coding literature [12, Fig. 7b], reproduced here in Figure 10. We consider two different messages, $M_1, M_2 \sim \text{i.i.d. Ber}(1/2)$, provided as input to the system. Edges along which information about M_1 flows are colored in blue, while edges along which information about M_2 flows are colored in orange. The reader may identify these using Definition 4 and the transmission on each edge shown in Figure 10.

An important feature to observe is that when C_2 mixes information by computing the XOR of M_1 and M_2 , we see information about M_1 spontaneously beginning to flow on (B_2, B_3) and similarly, information about M_2 beginning to flow on (A_2, A_3) . This is expected, since M_2 is relevant for decoding M_1 at this stage, and indeed, it is exactly this idea which is used to decode M_1 at B_4 . All of this is true, despite the fact that $M_1 \oplus M_2$ is independent of M_1 and M_2 individually. This is once again, a prime example of synergy in action.

Applying the information path algorithm for the message M_1 at A_4 will reveal two paths: the “upper path”

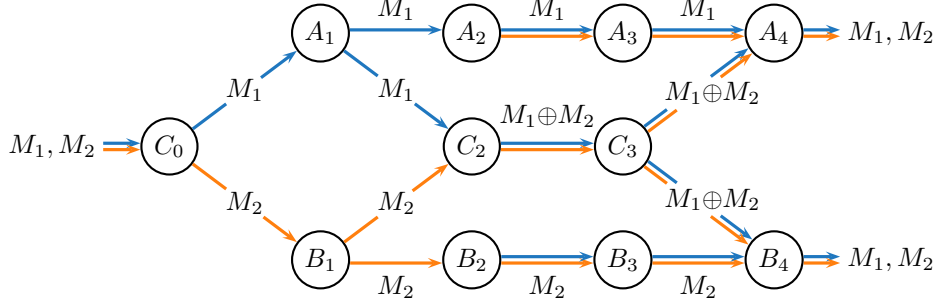


Figure 10: A depiction of the butterfly network discussed in Section 6.1. There are two messages, M_1 and M_2 , each with its own information flow. All edges with M_1 -information flow are shown in blue and all edges with M_2 -information flow are shown in orange. After time $t = 2$, all edges shown have both M_1 - and M_2 -information flow. Once the system computes $M_1 \oplus M_2$, edges transmitting M_1 have information flow about both M_1 and M_2 , since M_2 can now be decoded from $M_1 \oplus M_2$ and M_1 . Furthermore, observe the M_1 - and M_2 -information paths in this system. In particular, there are two possible M_1 -information paths to A_4 , but only one possible M_2 -information path, which flows through the middle link. The same applies to the M_1 -information path to B_4 . This may suggest the importance of the middle link in enabling this computation.

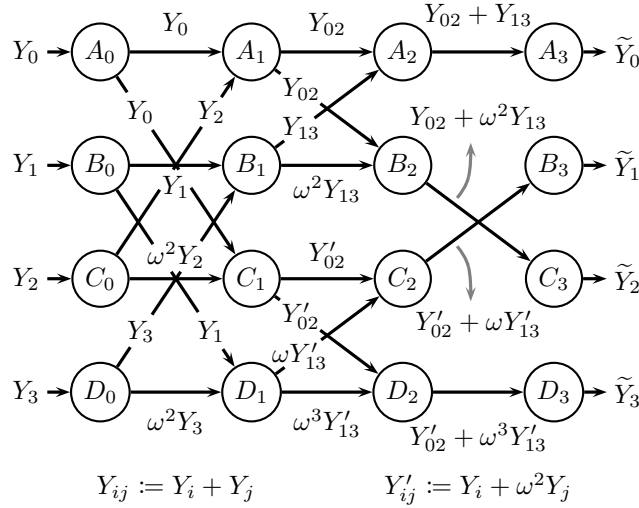


Figure 11: The computational system of the 4-point Fast Fourier Transform. For brevity, we have set $\omega := e^{-j\frac{2\pi}{4}}$.

$(C_0, A_1, A_2, A_3, A_4)$, and the “middle path” $(C_0, A_1, C_2, C_3, A_4)$. However, applying the information path algorithm for the message M_2 at A_4 reveals that M_2 exclusively used the “middle path”, $(C_0, B_1, C_2, C_3, A_4)$, to arrive at A_4 from the input nodes.

6.2 The Fast Fourier Transform

The Fast Fourier Transform (FFT) is a well-known computational network that provides an intuitive setting for examining information flow. In general, the N -point FFT is an implementation of the N -point Discrete Fourier Transform (DFT), given by

$$\tilde{Y}_k = \sum_{i=0}^{N-1} Y_i e^{-j\frac{2\pi k}{N}i}, \quad k \in \{0, 1, \dots, N-1\}. \quad (67)$$

The DFT is a basis transformation of a discrete-time signal Y , which is usually assumed to be periodic with period N . The N -point DFT represents such a signal in the complex-exponential Fourier basis, yielding the Fourier coefficients \tilde{Y} . We consider a simple 4-point DFT, i.e. $N = 4$. The FFT implements this transform

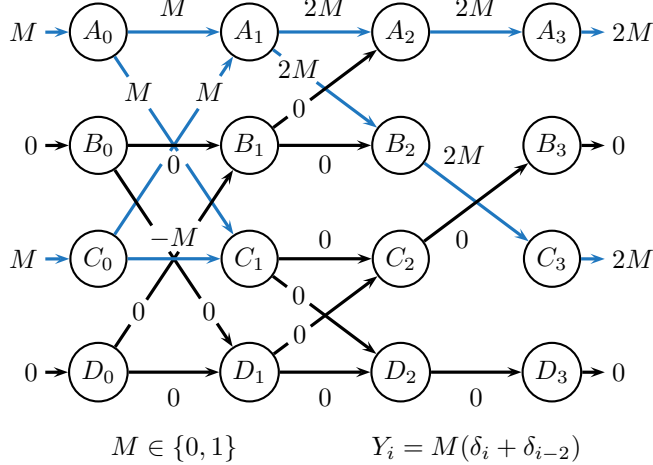


Figure 12: An example of information flow in the 4-point FFT, when the message determines which of two signals is supplied to the system: $Y = [0, 0, 0, 0]$ or $Y = [1, 0, 1, 0]$. Observe that, since M is encoded in the even part of Y , only the “even component” of the FFT network is active. Furthermore, only the DC component, \tilde{Y}_0 and the first harmonic, \tilde{Y}_2 are active, as we would expect based on the two input signals.

using the computational system shown in Figure 11. We refer the reader to [44, Ch. 9] for details. For notational convenience, we have set $\omega = e^{-j\frac{2\pi}{N}} = e^{-j\frac{2\pi}{4}}$.

We use this example to demonstrate how the definition of the message is important in determining information flow. First, suppose the message is one of two signals: $Y = [0, 0, 0, 0]$, or $Y = [1, 0, 1, 0]$. This can be written as $M \in \{0, 1\}$ and $Y_i = M(\delta_i + \delta_{i-2})$, where $\delta_i = \mathbb{I}\{i = 0\}$ is the Kronecker Delta function, and we assume $M \sim \text{Ber}(1/2)$. The full computational system, along with the random variables computed on all edges, is shown in Figure 12. The edges that have M -information flow are highlighted in blue. Since M is encoded into the *even* part of Y (observe that $Y_i = Y_{-i} \forall M$), we notice that only the “even component” of the FFT system (corresponding to the 2-point FFT on the even indices of Y) is active [44, Sec. 9.3]. Furthermore, only \tilde{Y}_0 , the DC component, and \tilde{Y}_2 , the first harmonic, show variation with M at the output, as we would expect based on the two input signals.

As a second example, consider the case shown in Figure 13. Here, the message is again one of two signals: $Y = [1, 1, 1, 1]$, or $Y = [1, 1/\omega, 1/\omega^2, 1/\omega^3]$. These signals can be jointly expressed in terms of the binary message random variable $M \sim \text{Ber}(1/2)$ as $Y_i = 1/\omega^{iM}$. The two signals are flat in their magnitude spectra and differ only in their phase, creating δ -functions in the Fourier domain that are frequency-shifted with respect to one another: $\tilde{Y}_k = \delta_{k-M}$. Once again, the edges in the network that carry M -information flow are demarcated in blue. A detailed derivation of the values of the transmissions in the computational system can be found in Appendix F.1.

These two examples make it clear that, based on how the message is defined, the M -information paths in the system can be very different. Indeed, if the message were as general as possible, by placing a probability distribution over all possible values of Y in \mathbb{R}^4 , we know that *all* edges in the computational system would have M -information flow. However, selectively restricting M to just a few signals helps reveal some kind of structure within the FFT network.

Another feature that can be observed in these examples is how the output of the computational system can be a function of the message. Although only very simple functions of the message have been shown at the outputs here, the FFT demonstrates that, in principle, more complex functions of the message may also be generated.

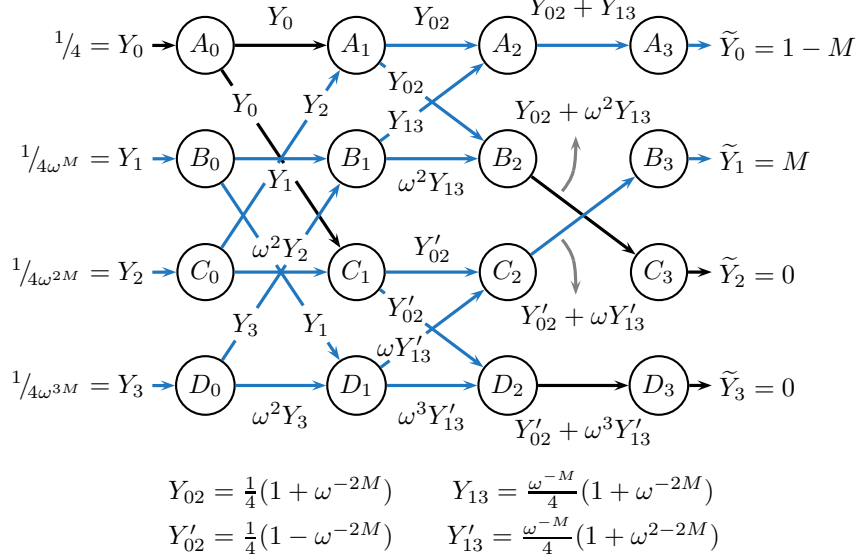


Figure 13: Another example of information flow in the 4-point FFT, when the message determines which of two signals is supplied to the system: $Y = [1, 1, 1, 1]$ or $Y = [1, 1/\omega, 1/\omega^2, 1/\omega^3]$. The M -information paths are different from those in Figure 12, showing how the choice of the message can have a strong impact on the flows within the same computational system.

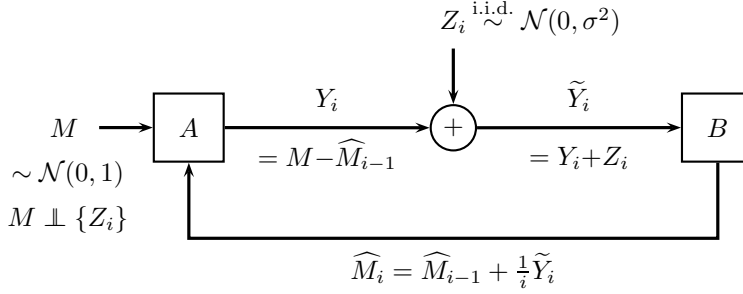


Figure 14: A communication system depicting the Schalkwijk and Kailath scheme. Alice, represented by node A , communicates a message M to Bob, represented by node B , in the presence of a noisy feedforward channel and a noiseless feedback channel. In the i^{th} iteration, Alice transmits the error in Bob's most recent estimate of the message, Y_i , but her transmission is corrupted by the noise Z_i . Bob updates and transmits his estimate, \widehat{M}_i , which reach Alice noiselessly.

6.3 The Schalkwijk and Kailath Scheme

The Schalkwijk and Kailath scheme [71] is an efficient strategy for communicating a message in the presence of a noisy feedforward channel and a noiseless feedback channel. We have previously used this scheme as a counterexample [41], to show that comparing Granger causal influences in forward and backward directions can lead to erroneous inferences on the direction in which the message is being sent in this feedback system. We first provide a brief overview of the scheme, then recapitulate our previous result, and finally demonstrate what the information flow framework developed in this paper has to offer in the case of this example.

Consider the communication system depicted in Figure 14, which shows the schematic of a simplified version of the Schalkwijk and Kailath scheme. For convenience, let us denote the transmitter, A , and receiver, B , by Alice and Bob respectively. Alice is attempting to communicate a message M to Bob over an additive Gaussian channel, but in the presence of noiseless feedback. Alice starts by transmitting the message $Y_1 = M$ to Bob, over the noisy feedforward channel. Bob receives a corrupted version of M , given by $\widetilde{Y}_1 = Y_1 + Z_1$, and computes an estimate \widehat{M}_1 . He sends this estimate back to Alice over the noiseless feedback channel. In the iterations that follow, Alice computes the error in Bob's most recent estimate, $Y_i = M - \widehat{M}_{i-1}$, and

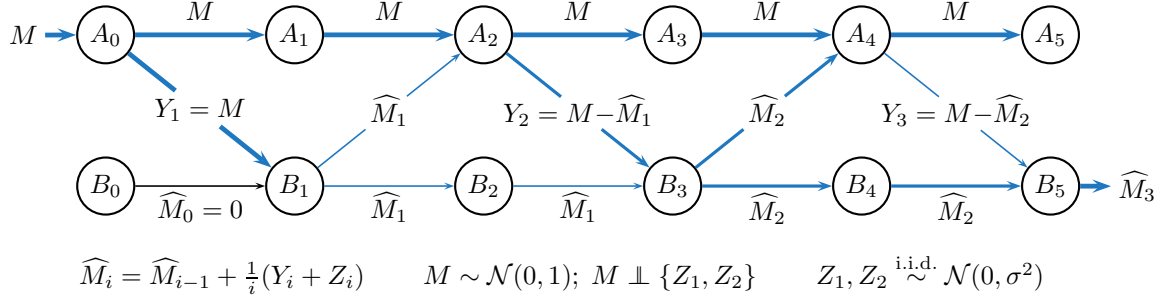


Figure 15: A computational system describing the first few iterations of the Schalkwijk and Kailath scheme. Almost every edge shown here has M -information flow. However, the *quantity* of M -information flow (shown using line thickness) reveals the asymmetry between Alice and Bob: Alice has the message to begin with, and her transmissions have a larger volume of M -information flow. In contrast, Bob’s initial transmissions are poor estimates and have small volumes of M -information flow, but they get better over a few iterations, and eventually come close to the true message. Furthermore, we also reveal an asymmetry between Alice and Bob using the concept of derived information: each of Bob’s transmissions is M -derived from Alice’s previous transmissions, whereas Alice’s transmissions are *not* M -derived from Bob’s previous transmissions. Both these facts point towards the idea that Alice is slowly sending information about M to Bob.

sends this to Bob over the noisy feedforward channel. Meanwhile, Bob updates his estimate based on Alice’s noisy transmissions $\tilde{Y}_i = Y_i + Z_i$, using the following rule:

$$\widehat{M}_i = \widehat{M}_{i-1} + \frac{1}{i} \tilde{Y}_i \quad (68)$$

It can be shown that this rule implies

$$\widehat{M}_i = M + \frac{1}{i} \sum_{j=1}^i Z_j \quad (69)$$

Thus, this strategy ensures that Bob’s estimate \widehat{M}_i converges to M in mean squared sense [41].

Intuitively, one might expect that, since the message M is being transmitted in the forward direction, the Granger Causal influence from Alice to Bob is greater than that from Bob to Alice. However, our earlier result [41] showed that, in fact, the opposite is true. In other words, even though the message is being communicated from Alice to Bob, the Granger Causal influence from Bob to Alice is greater; in fact, the Granger Causal index from Bob to Alice is *infinite*. The reason for this is that, while Alice’s past transmissions do not perfectly predict Bob’s transmissions (due to the presence of noise in the feedforward link), Bob’s past transmissions *perfectly* predict Alice’s transmissions (since the latter are a simple function of the former). Therefore, the Granger Causal index from Alice to Bob, which measures the relative predictive gain of including Alice’s past transmissions in the autoregression for Bob’s transmissions, remains finite; while the Granger Causal index from Bob to Alice becomes infinite.

Our earlier paper on this subject [41] concluded that the direction of greater Granger Causal influence could be opposite to the “direction of information flow” in the Schalkwijk and Kailath scheme. There, “information flow” was being used purely in an intuitive sense, to mean the direction in which the message was being communicated in that system. The intent of our previous paper was to explain that it is not always possible to interpret a larger Granger causal influence in a certain direction to mean that a specific message is being communicated in that direction. In contrast, this paper presents a refined theoretical framework that defines information flow about a message M for a specific *edge* in a computational system. Now, we no longer speak of *one specific direction* in which information flows; rather, we describe *which edges* carry information about the message in their transmissions *at each point in time*. This leads to a more nuanced understanding of information flow in the Schalkwijk-Kailath setting.

Before we can analyze the M -information flows in the Schalkwijk-Kailath scheme, we need to fit the scheme within the computational system framework. Figure 15 shows the time-unrolled computational system cor-

responding to two feedforward and feedback iterations of the simplified Schalkwijk-Kailath scheme described before. In order to translate the communication system into our computational system model while remaining consistent with our earlier work [41], we have merged the process of noise addition with the receiver, i.e., Bob. This exposes the edges with Alice’s and Bob’s *transmissions*, making them observable, as was assumed in our previous paper [41]. This is also consistent with what *would have been* observable if A and B were neurons (or neural populations) whose outputs a neuroscientist were to measure.²⁶ Note that one full iteration of the Schalkwijk-Kailath scheme takes two time steps in this model, so the iteration index i advances once for every two time steps t . Also, note that this does *not* make \tilde{Y} or Z “hidden nodes”, since the function computed at B_t can be defined purely in terms of its inputs, (Y_i, \widehat{M}_{i-1}) , and its intrinsic random variable, $W(B_t)$ (which absorbs Z_i), as follows:

$$f_{B_{2i-1}}(Y_i, \widehat{M}_{i-1}, W(B_{2i-1})) = \widehat{M}_{i-1} + \frac{1}{2}(Y_i + W(B_{2i-1})) \quad (70)$$

where $W(B_{2i-1}) = Z_i$ takes the role of the noise in the communication system. Also, to understand the time index for node B , note that in the first step of iteration i , Alice transmits to Bob, i.e., node A_{2i-2} transmits to B_{2i-1} (see Figure 15).

Now, we first show that all edges depicted in blue in Figure 15 carry M -information flow, based on Definition 4. Specifically, both Alice’s feedforward transmissions *and* Bob’s feedback transmissions have M -information flow. This should not be surprising for the following intuitive reasons: Alice’s transmissions convey information about M which Bob uses to improve his estimate; meanwhile, Bob’s transmissions are estimates of M , and therefore must depend on M . In fact, we can take this intuitive argument further: suppose we were to *quantify* M -information flow by using the following natural extension of our definition,

$$\mathcal{F}_M(E_t) := \max_{\mathcal{E}'_t \subseteq \mathcal{E}_t} I(M; X(E_t) \mid X(\mathcal{E}'_t)). \quad (71)$$

Noting that Definition 4 only specified *whether or not* a given edge E_t had information flow, all that we have now done is to take the maximum over the subsets of edges used to discover M -information flow in that definition. This quantification is fully consistent with our definition of M -information flow, since it goes to zero if and only if the M -information flow on an edge goes to zero. Now, using this quantitative notion of information flow, we can ask how the M -information flow on a given link—feedforward or feedback—varies with time. In particular, it should be intuitively clear that the M -information content in Bob’s transmissions, i.e. \widehat{M}_i , *increases* over time as his estimate improves. This is depicted as an increase in the thickness of the edges carrying Bob’s transmissions with time. Meanwhile, the information content in Alice’s transmissions *decreases* with time. To understand why the latter is true, note that, after the first iteration, Alice’s transmission represents the noise in Bob’s estimate. Therefore, just as in Counterexample 1, when conditioned on Bob’s estimate, Alice’s transmissions depend on the message. At the initial iterations, when Bob’s estimate is poor, we must have that $I(M; \widehat{M}_i)$ is very small. Hence, we see that:

$$I(M; Y_i \mid \widehat{M}_{i-1}) = I(M; M - \widehat{M}_{i-1} \mid \widehat{M}_{i-1}) \quad (72)$$

$$\stackrel{(a)}{=} H(M \mid \widehat{M}_{i-1}) + H(M \mid M - \widehat{M}_{i-1}, \widehat{M}_{i-1}) \quad (73)$$

$$= H(M \mid \widehat{M}_{i-1}) \quad (74)$$

$$= H(M) - I(M; \widehat{M}_{i-1}) \stackrel{(b)}{\approx} H(M), \quad (75)$$

where in (a), the second term is zero because M is a constant when given \widehat{M}_{i-1} and $M - \widehat{M}_{i-1}$; while in (b), $I(M; \widehat{M}_{i-1})$ is assumed to be approximately zero when Bob’s estimate is poor (as we might expect if

²⁶From a wireless communication system perspective, as well, it is more reasonable to assume noise to be a part of the receiver’s node, since the additive noise in a signal is usually considered to be the result of thermal noise in the receiver’s circuitry.

the noise is large, for instance). Hence, for the first few iterations, the quantified M -information flow of Alice’s transmissions is close to $H(M)$, which is as large as the flow can get. However, as Bob’s estimate improves, $I(M; \widehat{M}_{i-1})$ becomes closer to $H(M)$, and therefore $I(M; Y_i | \widehat{M}_{i-1})$ becomes close to zero. At the same time, $I(M; Y_i)$ is equal to zero, since Y_i carries only information about the noise in \widehat{M}_{i-1} (after the first iteration), which is independent of M . Thus, the quantified M -information flow of Alice’s transmissions decreases over time. Correspondingly, this is depicted using edges whose thickness decreases over time in Figure 15.

Quantifying the M -information flows of the feedforward and feedback links thus reveals an asymmetry between Alice and Bob that strongly suggests that the message is being transmitted from Alice to Bob. However, we can get a more nuanced understanding of information flow in this system by asking whether Bob’s transmissions are *derived* from Alice’s, or vice versa. First, consider whether Bob’s transmissions are derived M -information of Alice’s previous transmissions: this can be expressed in terms of the Markov chain $M \rightarrow [M, M - \widehat{M}_1] \rightarrow \widehat{M}_2$. Observe that this Markov chain holds trivially:

$$I(M; \widehat{M}_2 | M, M - \widehat{M}_1) = 0. \tag{76}$$

However, if we consider whether Alice’s transmissions are derived M -information of Bob’s past transmissions, it can be shown that $M \rightarrow [\widehat{M}_1, \widehat{M}_2] \rightarrow (M - \widehat{M}_2)$ is not a valid Markov chain (see Appendix F.2 for a detailed derivation). Hence, we see that Bob’s transmissions are derived M -information of all of Alice’s past transmissions, however, Alice’s transmissions are *not* derived M -information of all of Bob’s past transmissions. In conjunction with the fact that the volume of M -information flow in Alice’s transmissions slowly decreases from $H(M)$ with time, while the volume of M -information flow in Bob’s transmissions slowly increases to $H(M)$ with time, this suggests that Alice has some information about the message M that Bob slowly receives from Alice.

This example shows how a measure that quantifies information flow, along with derived information, can be used to understand some finer computational structure present within the computational system. In general, however, care needs to be exercised in applying derived M -information: one must choose what Markov condition to check in a principled manner. In the specific case of the Schalkwijk-Kailath example, we had the advantage of being in a two-node setting, where the derived information expressions we examined had clear interpretations. It may be that analyzing information flow first, to understand which variables transmit information about M to one another, can help guide the choice of variables to examine when applying derived M -information.

6.4 A Message Defined at the Output of a System

We now describe an example where the message is defined at the *output* of a computational system, instead of at the input. Although Definition 3c defines the message to be a random variable available at the input nodes, it is also possible to define the message at the output of the computational system. In this scenario, the input nodes are no longer well-defined as per Definition 3c. Instead, we would define *output nodes* in the same manner.²⁷

Consider the computational system shown in Figure 16. The system on the right executes the function depicted by the boolean circuit shown on the left. $Y \in \{0, 1\}$ is an external parameter, which is taken to be a fixed constant. When $Y = 1$, the AND gate at the top is activated while the AND gate at the bottom

²⁷Note, however, that the corresponding “opposite” of Theorem 7 (wherein the places of “input” and “output” nodes are switched) does not hold in this case. That is, it is *not true* that if at some previous time instant, an “*input*” node’s outgoing transmissions depend on the message, then there is an information path connecting that input node to the aforementioned output nodes. The reason this fails is that there could be a “source” node at an even earlier time instant, which provides information about M to *both* the input node under consideration, and the output nodes, via two separate, diverging paths. Therefore, there may be no path from the input node to the output nodes.

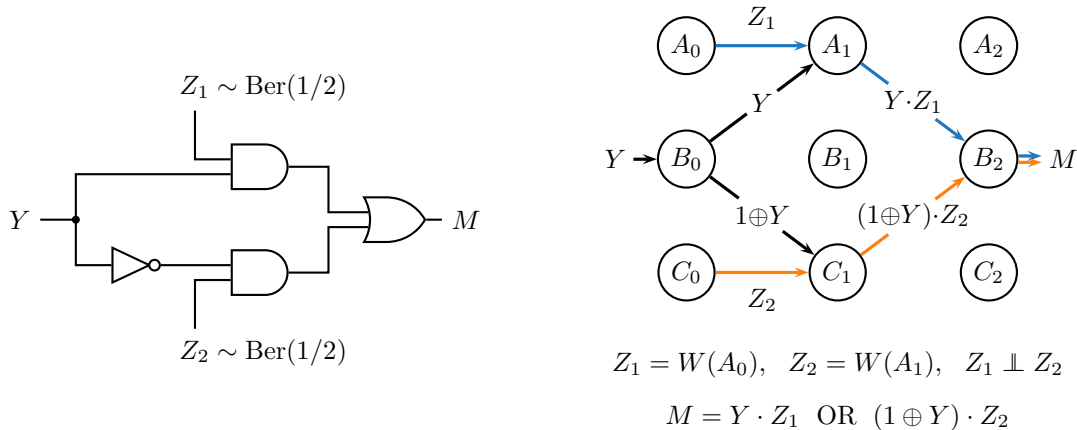


Figure 16: A boolean circuit demonstrating a message defined at the output of the computational system. Note that “ \oplus ” refers to bitwise-XOR, “OR” refers to bitwise-OR, and “ \cdot ” refers to bitwise-AND. We see that information paths may lead from an internal node, that generates an intrinsic random variable, to the output node. Furthermore, this path may change with the “external parameters” of the system.

is deactivated, so the message depends only on Z_1 . In this case, only the edges shown in blue have M -information flow. On the other hand, when $Y = 0$, the opposite happens, and the message M depends only on Z_2 . Now, only edges shown in orange have M -information flow. If Y was not a deterministic external parameter, but a random variable itself, then all edges shown in the figure would have M -information flow, since M would depend on all their values.

So, we see that when the message is defined at the output, the “origin” of the message may be from within the computation system itself, in the form of one or more intrinsically generated random variables: here, either $Z_1 = W(A_0)$ or $Z_2 = W(C_0)$. The notion of information flow and information paths can thus help us identify where the message originates within the computational system.

Furthermore, just as information paths can change depending upon how the message is defined (as in Section 6.2), information paths may also change depending on external parameters: inputs such as Y that are fed into the computational system, which are not part of the message. These inputs essentially shape the nature of the computation being performed, and so, naturally, they can affect information paths.

7 Discussion

This paper presented a theoretical framework for defining and studying information flow of a specific message in a computational system. The core contribution of our paper was a definition for information flow that is concretely grounded in the *computational task* and intimately tied to a *specific message*. This relied on another important contribution: the development of an underlying computational model, which enables the interpretation of statistical analyses. After providing a clearly-defined model for a computational system, we presented several candidate definitions for information flow along with counterexamples and showed that our definition, which is based on positivity of a conditional mutual information expression, satisfies several intuitive properties, whereas other candidate definitions do not. We then examined these properties in detail and showed, in particular, that our definition naturally leads to the existence of “information paths”. We also discussed how information flow can be inferred through conditional independence testing, and provided an algorithm for recovering the information paths in a given system. Finally, we studied some canonical examples of computational systems from different contexts, and showed that our definition of information flow is intuitive in each case.

We proceed to discuss several important assumptions and simplifications in our model. We also discuss existing literature related to estimation of causal influence in neuroscience, and how our computational system model leads to a significantly different measure of information flow. Similarly, we discuss how our framework is very different from the field of Probabilistic Graphical Models.

7.1 Neuroscientific Concerns

7.1.1 Observing edges vs. nodes

The observation model stated in Section 5.1 makes a crucial assumption, namely, that transmissions on each *edge* can be observed. In neuroscientific experiments, however, we often record activity from single neurons (as in the case of electrophysiological recordings), or aggregate activity from groups of neurons (as with Local Field Potentials measured in Electrocorticography and Electroencephalography). These neurons, or groups of neurons, are considered to be nodes communicating to one another in a network. It may not be known which nodes are connected to which other nodes, let alone the recipient of each transmission at every time instant. This is a marked departure from our assumption that transmissions on edges can be observed. To some extent, it is possible to incorporate a “node-centric” model within our computational system by assuming that all nodes broadcast their transmissions. However, that still leaves unanswered the question of which nodes actually “hear” another’s transmissions. A possible resolution to that question might arise from an understanding of *receiver response*. That is, we consider a revised model in which an edge exists if a receiving neuron *uses* the information transmitted by some neuron at the previous time instant. This issue is beyond the scope of the current work, and will be addressed in subsequent studies.

We also note that, although tools based on Granger causality implicitly assume that nodes are measured and not edges, they do not resolve the issue of which node is “talking” to which other node. For example, if two different nodes A_1 and B_1 communicate the same information to a third node, C_2 , any regression based analysis will assign a weight of one-half to each of A_1 and B_1 . However, the true function, f_{C_2} , may be using only the information coming from A_1 , or only the information coming from B_1 , or using the two in some other unequal proportion. Such cases may only be identifiable through an interventional approach.

7.1.2 Observing memories

Another important assumption in the observation model is that *memories* of nodes are observed as transmissions on self-edges. If these transmissions are implemented in the form of some internal state at each node, then they might be difficult to observe in practice.²⁸ It remains to be fully understood whether one can compensate for not observing memories in some manner, e.g., by assuming that the memory of a node is the full history of its transmissions and receptions. While this means that intrinsically generated random variables that are *not* propagated to other nodes will never be observed, it could be argued that such variables could have no impact on the system (save for acting as “computational noise”). So perhaps it suffices to observe only transmissions between *different* nodes (and not self-edges). Further work is required to understand what ramifications such an assumption has on identifying information flows and information paths.

Conversely, our work may suggest to neuroscientists that inferences about information flow are more reliably obtained if one can measure transmissions on edges in the graph, rather than transmissions of nodes. This may call for newer imaging modalities, or new uses of existing modalities, such as treating axons as targets

²⁸If every node represents a *group* of neurons, however it may just be that their internal state is represented in the form of *communication between these neurons*. In that case, perhaps observing their internal state is just a matter of having more spatially refined measurements.

for invasive recordings, perhaps at nodes of Ranvier. Further, perhaps if one wishes to observe memories, it is important to measure not only spikes, but also membrane voltages (e.g. using voltage-sensitive dyes [72] or, less directly, through measurements of changes in neurotransmitter concentrations outside a cell [73]).

7.1.3 Discretization of time

Yet another implicit assumption in our computational system model is that transmissions occur at discrete points in time. This assumption is justified for synchronous digital circuits used commonly today, or if the computational system of interest is a trained artificial neural network, for instance. However, this is not a perfect model of the brain, because neural spiking (among other processes), does not occur only at multiples of some fundamental unit of time. This issue might be partially mitigated by assuming that neural computation happens at a certain time scale, and by using a sufficiently high sampling rate so that Nyquist-rate-type arguments apply. This may not be possible in certain modalities (e.g. Calcium imaging and functional Magnetic Resonance Imaging) that are inherently slow, however, so it would be interesting to understand what inferences we are no longer capable of making. Alternatively, if the sampling rate is too high, it may be useful to consider windows within which to look for M -information flow. The exact implications of using such preprocessing methods will also need to be studied in greater detail, in future.

7.1.4 Message enters at $t = 0$

Another assumption in our framework is that the message enters the system at, and only at, time $t = 0$. This is essential, given the way we have defined input nodes: nodes at time $t = 0$, whose outputs depend on the message (and which have no other shared source of randomness). However, this assumption does not allow for a dynamically evolving stimulus, which is also common in neuroscientific experiments. Suppose we allow the message to enter the system at a later time instant, say at some node U_t , for $t > 0$, i.e., U_t may compute a function not just of its inputs, but also of M . Then, if we want the information path theorem to continue to hold, we must also add U_t to the set of input nodes.²⁹ Thus, if we see dependence at some other node $V_{t'}$, at a later time instant $t' > t$, the information paths leading to $V_{t'}$ may arise from the original input nodes *or* from U_t , or both. As we might intuitively expect, the more time points we allow the message to enter *at*, the more such information paths we will likely see, making the results of our analysis harder to interpret.

7.1.5 Experimental design and the message

An important aspect of our work is that it explicitly incorporates the message, which in neuroscientific experiments is often some information contained in the stimulus. This aids the neuroscientist in designing experiments, for example, in understanding what stimuli will help them make a certain inference about information flow. In particular, one needs to use at least two different stimuli in order to obtain any determination about information flow. While this is implicitly understood in neuroscience, as evidenced by comparisons with baselines, or by the use of permutation tests to scramble stimulus-trial correlations for a null model, our framework provides a more direct method for identifying and interpreting stimulus-related information flow.

²⁹We should also expect that any Local Markovity conditions at time t (see Proposition 6) that involve the node U_t will no longer hold.

7.2 The Difficulty of Estimation

A strategy for detecting edges that have M -information flow was presented in Section 5.2. In practice, however, there are several issues associated with employing such a strategy. These are discussed below.

Firstly, we currently assume that observations are noiseless (see Section 5.1, Assumption 3). It is unclear, exactly, to what extent noisy observations will impact the inference of information flow. In particular, it is worth understanding whether small amounts of observation noise can be tolerated if all edges with M -information flow have a sufficiently large “volume” of information (i.e., the corresponding mutual or conditional mutual information is sufficiently large). As was described intuitively in Section 5.2, if the information volume is large, then even under noisy conditions, we might expect the test statistic to clear the threshold, so the presence of M -information flow can still be detected consistently. But small volumes of information that aggregate over time—e.g. information about M “trickling” over time from one node to another—could still pose issues. Such M -information flow could go undetected, as has been shown to occur in different contexts [38], using different measures of flow. It is possible that Derived Information, in particular, is hard to infer in the presence of noise. This could make the task of detecting the presence of a hidden node difficult (consider the case of a “trickling” hidden node), as well as that of identifying redundant links.

Secondly, detecting whether each edge at time t has information flow involves checking *all* subsets of \mathcal{E}_t . For N nodes and N^2 edges, this implies 2^{N^2} subsets of edges that need to be searched. This could be seen as being prohibitively difficult for $N^2 \geq 30$, or for N greater than about 5 or 6 nodes. However, in reality, graphs in neuroscience are often known to be edge-sparse [74–76]. For example, in the brain, a well-established 11-node network is the reward network [76]. Most nodes in this network typically have just one incoming and one outgoing connection. The two most important nodes have five incoming edges each, with two and four outgoing edges respectively. Further, it is known which connections are inhibitory and which are excitatory, which could further help with testing for information flow. A fully connected network would have had 121 edges, but the underlying connectivity of the circuit only allows for a total of 17 edges in this network. So in reality, anatomical priors help reduce the number of edges to well within the range of what is computable. Nevertheless, it remains of interest to find methods by which nodes and/or edges can be excluded from the search, and this could be another topic for further research.

Another statistical issue that crops up when attempting to simultaneously perform several conditional independence tests is the problem of *multiple comparisons* [77]. Simply put, when performing a large number of independent hypotheses tests, say N , at some fixed false alarm rate α , on average, we should expect αN of these tests to erroneously reject the null. In the context of information flow, we might wish to set the null hypothesis to be the absence of M -information flow on a given edge. Then, to test for M -information flow on this edge, we need to perform a large number of conditional independence tests—call this number N —at some false alarm rate α . These tests are, in fact, *not independent* of one another; nevertheless, very loosely put, if we choose a false alarm rate $\alpha \approx 1/N$, we may find that the probability of *at least one* false alarm is too high. This would make us erroneously infer that this particular edge has M -information flow; moreover, since this argument applies to any edge, if α is not chosen conservatively enough, we may erroneously infer that *all* edges have M -information flow. This multiple hypothesis testing problem is better posed as a “Global Null test” (e.g., see [78]), wherein the global null is the hypothesis that *all* of the conditional independence tests are individually null (i.e., that there is *no* M -information flow on the given edge), and the global alternative is the hypothesis that *at least one* of the conditional independence tests is non-null (i.e., that there *is* M -information flow on the given edge). As mentioned before, however, the conditional independence tests dictated by Definition 4 are, in general, *dependent* on one another. Furthermore, it might not be easy to describe the manner of dependence, so when choosing methods that control the family-wise error rate, it is essential to choose those that work under arbitrary dependence. A simple example of such a test is the well-known Bonferroni correction, which uses a level $\alpha' = \alpha/N$ for each test

(where α is the desired false alarm rate for the overall global null test); but we may find that such methods have insufficient statistical power. A potential solution to this problem might involve combining multiple global null tests in some meaningful way: for example, one could imagine designing a procedure that controls the False Discovery Rate³⁰ [79] on the *identification of edges* with M -information flow.³¹ Another approach might be to find ways of directly testing information *paths*, wherein the hypothesis tested would be that a certain M -information *path* exists in the system, rather than requiring every *edge* with M -information flow be identified first. All of these ideas are potential avenues for future work.

7.3 The Limitations of Granger Causality and Related Tools

Mapping directed functional connectivity and information flow in the brain has been a hot topic for several years, as evidenced by the large body of work in this direction [20–22]. Approaches for statistically mapping functional connectivity often rely on variations of Granger Causality [24] and, more recently, Directed Information [26–28], which we here collectively refer to as “Granger Causality-based tools”. These approaches lack a systematic framework that ties the statistical analysis to the underlying computation, however, and the interpretations drawn from their use have often been questioned [31, 33, 34, 38–41].

In particular, a crucial difference between our approach and that of Granger Causality-based tools is that the latter do not have an explicit description of the message. Instead, they provide mechanisms to condense a pair of time series into a single statistic. There are no concrete models that can be used to interpret what this statistic means for the flow of *information* about the message. Furthermore, if one is interested in the information flow of multiple messages, Granger Causality-based tools do not provide an immediate solution. This is why a tool that ties information flow directly with a message is of great interest to practitioners.

The absence of an underlying computational framework with well-defined assumptions inherently makes it very hard to draw sound inferences through the application of Granger Causality-based tools. A striking example of this is a recent result of ours [41] that shows, using a feedback communication system, that the direction of greater Granger-causal influence can be opposite to the direction in which the message is communicated, even in the absence of hidden nodes and measurement noise. The time-unrolled graph framework presented here has been specifically designed to address this issue, and present a clear understanding of information flow, even in the presence of feedback. The example given in Section 6.3 demonstrates a potential resolution to this issue.

Granger Causality was originally developed for the study of time-series that occur only once, such as in economics [23]. An artifact of this development is that it was not designed to incorporate multiple trials of the same process. Instead, it assumes stationarity to help estimate parameters of the random variables that control the process. In the neuroscientific context, stationarity is often a very poor assumption, since the segment of time-series data corresponding to each trial may be short, and often sees some kind of stimulus presentation. Naturally, presentation of the stimulus changes the underlying parameters of the time-series and destroys stationarity; indeed, this is the quintessential aspect of the experiment. Thus, in order to understand processing in such stimulus-driven tasks, one needs to be able to infer time-dependent information flows from data. While information-theoretic extensions of Granger Causality such as Transfer Entropy and Directed Information do not assume stationarity, they nevertheless fail to provide a dynamically evolving picture of information flow.

Lastly, it is unclear whether the directional influences estimated using Granger Causality-based tools have any correspondence with the rigorous notion of information flow we have derived here, under special assumptions,

³⁰These methods control the expected proportion of false discoveries, i.e., the proportion of null hypotheses that are falsely rejected.

³¹Care is needed when doing this, however, since tests for M -information flow on different edges at the same time instant are *also* dependent on one another.

e.g., Gaussianity. This is a promising future direction as well, since it is important to understand in which situations these methods recover meaningful flows of information, and in which cases we must be careful with interpretation.

7.4 Probabilistic Graphical Models and Pearl’s Causality

There is one important difference that distinguishes our work from the perspective adopted in the field of probabilistic graphical models (PGMs) [80], and the representations therein. In our framework, nodes represent computational units, whereas in PGMs, nodes represent the random variables themselves, and edges capture the conditional independence relationships between these variables. While it might be possible to construct a PGM that is equivalent to our computational model, this would likely eliminate any intuitive structure captured by the computational graph.

It remains to be understood whether and how Pearl’s notions of causality [13] can be seamlessly merged with the understanding of information flow developed here. We expect that some formal application of causality will be needed in going from an edge-centric model (as presented here) to a more node-centric one (discussed in Section 7.1), in order to identify which transmissions influenced a given node’s output.

There are several works in the literature that discuss measures of information flow in probabilistic graphical models [81, 82], but they are heavily inspired by causality and largely center around an interventionist approach. In contrast, our definition of information flow is based on a computational system model that translates more readily to neuroscience, and we assume that the experimentalist is restricted to making observations.

7.5 Future Directions for Theoretical Development

A natural question that arises from this paper is: how can our definition of information flow on an edge be extended to a more generic information *measure*, which also quantifies the volume of flow? Finding such a measure will involve aggregating the conditional mutual information for each subset of edges into a single value (one example of such a measure was provided in Section 6.3, though it was not developed from first-principles). It is as yet unclear how this might be achieved, while still gelling well with our intuition of what this information flow volume ought to be. We believe that the right approach is to start by designating a set of properties that we would like information flow volumes to satisfy, and then to propose a measure through the use of representative examples and counterexamples.

A second direction that emerges is related to Partial Information Decomposition (PID) [47–49], which was discussed earlier in Section 3.5. M -information flow is very closely related to the PID: while Candidate Definition 1 checks for positivity of mutual information between M and $X(E_t)$, and hence implying the presence of unique and/or redundant information, our definition also detects the presence of purely synergistic information. Since our definition is closely tied to computation and is strongly motivated through the goal of finding unbroken information paths, the close relationship between PID and our definition suggests that PID might be the right toolset for obtaining a more fine-grained understanding of information flow, as well as computation. In particular, it would be useful to know how the understanding of computation is enhanced through a PID analysis, which describes the unique, redundant and synergistic components of the message in different nodes’ transmissions. Finally, we note that the PID could also help inform the discussion on a definition for information volume. Providing a useful definition of information volume based on current definitions of unique, redundant and synergistic information, and asking whether the problem of information flow can inform the PID literature, will also be the subject of future research.

A third direction has to do with alternate definitions of information flow: the properties we stated in

this paper are not sufficient to uniquely specify our definition of information flow. For example, the all-zero function as well as the all-ones function satisfy the Broken Telephone property, although they are not particularly useful definitions of information flow. Thus, it would be useful to understand what other properties we should impose so as to arrive at a unique definition of information flow. As a crude and preliminary example, we demonstrate how this might be done in Appendix E.

7.6 Concluding Remarks

We conclude by describing some of our general impressions in working on the theoretical development presented in this paper. As such, these points merely highlight some of our opinions on how theory—and more specifically, information theory—may be applied in neuroscience.

As mentioned in the introduction, we drew inspiration from two papers that discuss how experimentalists understand systems in biology and neuroscience [15,16]. Both these works advocate for theory by arguing that we need new analytical tools, and that the accumulation of empirical knowledge alone does not constitute *understanding*. Lazebnik [15], in particular, mentions how terminology in biology tends to be vague and non-committal. We feel that an important reason for this is the absence of concrete underlying models, with clearly-stated assumptions. In other words, we think that theory and modeling can go a long way in providing a *language* that will enable well-grounded discussions. This language, in turn, arises through the development of theoretical models and formal definitions.

Another point made by both the aforementioned papers is that we should attempt to understand large computational systems by first examining smaller models, and models in which the ground truth is already known. This approach allows us to create new analytical tools that can be thoroughly vetted, so that the interpretations drawn from their use in experimental practice is unambiguous and undebated. We also believe that when trying to understand large computational systems, it is essential to start with toy models such as Counterexample 1. This philosophy of starting with toy models, and abstracting out meaningful ideas that hold more generally in large systems, is well-entrenched in the field of information theory, and can become a useful export in fields such as neuroscience.

Acknowledgments

We have many people to thank for extremely useful discussions. A non-exhaustive list follows: Mayank Bakshi, Marlene Behrmann, Todd Coleman, Uday Jagadisan, Haewon Jeong, Rob Kass, Gabe Schamburg, Tsachy Weissman. We also thank the anonymous reviewers whose comments improved our exposition substantially.

Praveen Venkatesh was supported, in part, by a Fellowship in Digital Health from the Center for Machine Learning and Health at Carnegie Mellon University. Pulkrit Grover was supported, in part, by an NSF CAREER Award.

References

- [1] Praveen Venkatesh, Sanghamitra Dutta, and Pulkrit Grover. How should we define information flow in neural circuits? In *2019 IEEE International Symposium on Information Theory (ISIT)*, pages 176–180, July 2019.

- [2] Jorge Almeida, Anat R. Fintzi, and Bradford Z. Mahon. Tool manipulation knowledge is retrieved by way of the ventral visual object processing pathway. *Cortex*, 49(9):2334–2344, 2013.
- [3] Andrea Brovelli, Mingzhou Ding, Anders Ledberg, Yonghong Chen, Richard Nakamura, and Steven L. Bressler. Beta oscillations in a large-scale sensorimotor cortical network: Directional influences revealed by Granger causality. *Proceedings of the National Academy of Sciences*, 101(26):9849–9854, 2004.
- [4] M. Bar, K. S. Kassam, A. S. Ghuman, J. Boshyan, A. M. Schmid, A. M. Dale, M. S. Hämäläinen, K. Marinkovic, D. L. Schacter, B. R. Rosen, and E. Halgren. Top-down facilitation of visual recognition. *Proceedings of the National Academy of Sciences*, 103(2):449–454, 2006.
- [5] Adam S. Greenberg, Timothy Verstynen, Yu-Chin Chiu, Steven Yantis, Walter Schneider, and Marlene Behrmann. Visuotopic cortical connectivity underlying attention revealed with white-matter tractography. *Journal of Neuroscience*, 32(8):2773–2782, 2012.
- [6] Constance Hammond, Hagai Bergman, and Peter Brown. Pathological synchronization in Parkinson’s disease: networks, models and treatments. *Trends in neurosciences*, 30(7):357–364, 2007.
- [7] Y Smith, MD Bevan, E Shink, and JP Bolam. Microcircuitry of the direct and indirect pathways of the basal ganglia. *Neuroscience*, 86(2):353–387, 1998.
- [8] Anthony A Grace. Gating of information flow within the limbic system and the pathophysiology of schizophrenia. *Brain Research Reviews*, 31(2-3):330–341, 2000.
- [9] Elodie Lalo, Stéphane Thobois, Andrew Sharott, Gustavo Polo, Patrick Mertens, Alek Pogosyan, and Peter Brown. Patterns of bidirectional communication between cortex and basal ganglia during movement in patients with Parkinson disease. *Journal of Neuroscience*, 28(12):3008–3016, 2008.
- [10] Jeffery Samuels and Nathan D. Zasler. *Event-Related Paradigms*, pages 1346–1347. Springer, 2018.
- [11] Clark David Thompson. *A Complexity Theory for VLSI*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, USA, 1980. AAI8100621.
- [12] R. Ahlswede, Ning Cai, S. Y. R. Li, and R. W. Yeung. Network information flow. *IEEE Transactions on Information Theory*, 46(4):1204–1216, July 2000.
- [13] Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2009.
- [14] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: Foundations and learning algorithms*. MIT press, 2017.
- [15] Yuri Lazebnik. Can a biologist fix a radio?—Or, what I learned while studying apoptosis. *Cancer cell*, 2(3):179–182, 2002.
- [16] Eric Jonas and Konrad Paul Kording. Could a neuroscientist understand a microprocessor? *PLoS computational biology*, 13(1):e1005268, 2017.
- [17] Pulkit Grover and Praveen Venkatesh. An information-theoretic view of EEG sensing. *Proceedings of the IEEE*, 105(2):367–384, Feb 2017.
- [18] P. Grover, J. A. Weldon, S. K. Kelly, P. Venkatesh, and H. Jeong. An information theoretic technique for harnessing attenuation of high spatial frequencies to design ultra-high-density EEG. In *2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 901–908, Sept 2015.

- [19] Praveen Venkatesh and Pulkit Grover. Lower bounds on the minimax risk for the source localization problem. In *2017 IEEE International Symposium on Information Theory (ISIT)*, pages 3080–3084, June 2017.
- [20] Karl J Friston. Functional and effective connectivity: a review. *Brain connectivity*, 1(1):13–36, 2011.
- [21] Karl Friston, Rosalyn Moran, and Anil K Seth. Analysing connectivity with granger causality and dynamic causal modelling. *Current opinion in neurobiology*, 23(2):172–178, 2013.
- [22] André M Bastos and Jan-Mathijs Schoffelen. A tutorial review of functional connectivity analysis methods and their interpretational pitfalls. *Frontiers in systems neuroscience*, 9:175, 2016.
- [23] Clive W J Granger. Investigating causal relations by econometric models and cross-spectral methods. *Econometrica: Journal of the Econometric Society*, pages 424–438, 1969.
- [24] Steven L. Bressler and Anil K. Seth. Wiener–Granger causality: A well established methodology. *NeuroImage*, 58(2):323–329, 2011.
- [25] James Massey. Causality, feedback and directed information. In *Proceedings of the International Symposium on Information Theory and its Applications (ISITA)*, pages 303–305, 1990.
- [26] Christopher J. Quinn, Todd P. Coleman, Negar Kiyavash, and Nicholas G. Hatsopoulos. Estimating the directed information to infer causal relationships in ensemble neural spike train recordings. *Journal of Computational Neuroscience*, 30(1):17–44, Feb 2011.
- [27] C. J. Quinn, N. Kiyavash, and T. P. Coleman. Directed information graphs. *IEEE Transactions on Information Theory*, 61(12):6887–6909, Dec 2015.
- [28] J. Jiao, H. H. Permuter, L. Zhao, Y. Kim, and T. Weissman. Universal estimation of directed information. *IEEE Transactions on Information Theory*, 59(10):6220–6242, Oct 2013.
- [29] Thomas Schreiber. Measuring information transfer. *Physical Review Letters*, 85:461–464, Jul 2000.
- [30] Luiz A. Baccalá and Koichi Sameshima. Partial directed coherence: a new concept in neural structure determination. *Biological Cybernetics*, 84(6):463–474, May 2001.
- [31] Olivier David, Isabelle Guillemain, Sandrine Saillet, Sebastien Reyt, Colin Deransart, Christoph Segebarth, and Antoine Depaulis. Identifying neural drivers with functional MRI: an electrophysiological validation. *PLoS biology*, 6(12):e315, 2008.
- [32] Alard Roebroeck, Elia Formisano, and Rainer Goebel. The identification of interacting networks in the brain using fMRI: model selection, causality and deconvolution. *Neuroimage*, 58(2):296–302, 2011.
- [33] Olivier David. fMRI connectivity, meaning and empiricism. comments on: Roebroeck et al. The identification of interacting networks in the brain using fMRI: model selection, causality and deconvolution. *Neuroimage*, 58(2):306–309, 2011.
- [34] Patrick A. Stokes and Patrick L. Purdon. A study of problems encountered in Granger causality analysis from a neuroscience perspective. *Proceedings of the National Academy of Sciences*, 114(34):E7063–E7072, 2017.
- [35] Lionel Barnett, Adam B. Barrett, and Anil K. Seth. Solved problems for Granger causality in neuroscience: A response to Stokes and Purdon. *NeuroImage*, 178:744–748, 2018.
- [36] Luca Faes, Sebastiano Stramaglia, and Daniele Marinazzo. On the interpretability and computational reliability of frequency-domain Granger causality. *F1000Research*, Sep 2017.

- [37] Patrick A Stokes and Patrick L Purdon. In reply to Faes et al. and Barnett et al. regarding “a study of problems encountered in granger causality analysis from a neuroscience perspective”. *arXiv:1709.10248 [stat.ME]*, 2017.
- [38] Jonas Andersson. Testing for Granger causality in the presence of measurement errors. *Economics Bulletin*, 2005.
- [39] Hariharan Nalatore, Mingzhou Ding, and Govindan Rangarajan. Mitigating the effects of measurement noise on Granger causality. *Physical Review E*, 75(3):031123, Mar 2007.
- [40] Mingming Gong, Kun Zhang, Bernhard Schölkopf, Dacheng Tao, and Philipp Geiger. Discovering temporal causal relations from subsampled data. In *Proceedings of The 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1898–1906. PMLR, Jul 2015.
- [41] Praveen Venkatesh and Pulkit Grover. Is the direction of greater Granger causal influence the same as the direction of information flow? In *53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 672–679, Sept 2015.
- [42] Claude Elwood Shannon. A mathematical theory of communication. *The Bell System Technical Journal*, 27(3):379–423, July 1948.
- [43] Karl J Friston, Lee Harrison, and Will Penny. Dynamic causal modelling. *Neuroimage*, 19(4):1273–1302, 2003.
- [44] Alan V. Oppenheim, John R. Buck, and Ronald W. Schafer. *Discrete-time signal processing*. Prentice Hall, Upper Saddle River, N.J., 2nd ed. edition, 1999.
- [45] Yury Polyanskiy and Yihong Wu. Lecture notes on information theory, August 2017.
- [46] Claude Elwood Shannon. Communication theory of secrecy systems. *Bell system technical journal*, 28(4):656–715, 1949.
- [47] Paul L Williams and Randall D Beer. Nonnegative decomposition of multivariate information. *arXiv:1004.2515 [cs.IT]*, 2010.
- [48] Malte Harder, Christoph Salge, and Daniel Polani. Bivariate measure of redundant information. *Phys. Rev. E*, 87:012130, Jan 2013.
- [49] Nils Bertschinger, Johannes Rauh, Eckehard Olbrich, Jürgen Jost, and Nihat Ay. Quantifying unique information. *Entropy*, 16(4):2161–2183, 2014.
- [50] Joseph T Lizier, Nils Bertschinger, Jürgen Jost, and Michael Wibral. Information decomposition of target effects from multi-source interactions: Perspectives on previous, current and future work. *Entropy*, 20(4):307, 2018.
- [51] Elad Schneidman, William Bialek, and Michael J Berry. Synergy, redundancy, and independence in population codes. *Journal of Neuroscience*, 23(37):11539–11553, 2003.
- [52] Peter E Latham and Sheila Nirenberg. Synergy, redundancy, and independence in population codes, revisited. *Journal of Neuroscience*, 25(21):5195–5206, 2005.
- [53] Nicholas M Timme and Christopher Lapish. A tutorial for information theory in neuroscience. *eNeuro*, 5(3), 2018.
- [54] Itay Gat and Naftali Tishby. Synergy and redundancy among brain cells of behaving monkeys. In *Advances in neural information processing systems*, pages 111–117, 1999.

- [55] Thomas M Cover and Joy A Thomas. *Elements of Information Theory*. John Wiley & Sons, 2012.
- [56] Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms*. The MIT Press, 3rd edition, 2009.
- [57] W. P. Bergsma. Testing conditional independence for continuous random variables. *EURANDOM report*, 2004(049), 2004.
- [58] Kun Zhang, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Kernel-based conditional independence test and application in causal discovery. In *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence, UAI'11*, pages 804–813, Arlington, Virginia, United States, 2011. AUAI Press.
- [59] Tzee-Ming Huang et al. Testing conditional independence using maximal nonlinear conditional correlation. *The Annals of Statistics*, 38(4):2047–2091, 2010.
- [60] Liangjun Su and Halbert White. A consistent characteristic function-based test for conditional independence. *Journal of Econometrics*, 141(2):807–834, 2007.
- [61] Meng Huang, Yixiao Sun, and Halbert White. A flexible nonparametric test for conditional independence. *Econometric Theory*, 32(6):1434–1482, 2016.
- [62] Rajat Sen, Karthikeyan Shanmugam, Himanshu Asnani, Arman Rahimzamani, and Sreeram Kannan. Mimic and classify: A meta-algorithm for conditional independence testing, 2018.
- [63] Rajen D. Shah and Jonas Peters. The hardness of conditional independence testing and the generalised covariance measure. *arXiv:1804.07203 [math.ST]*, Apr 2018.
- [64] Liam Paninski. Estimation of entropy and mutual information. *Neural computation*, 15(6):1191–1253, 2003.
- [65] Weihao Gao, Sreeram Kannan, Sewoong Oh, and Pramod Viswanath. Estimating mutual information for discrete-continuous mixtures. In *Advances in Neural Information Processing Systems*, pages 5986–5997, 2017.
- [66] Alexander Kraskov, Harald Stögbauer, and Peter Grassberger. Estimating mutual information. *Physical review E*, 69(6):066138, 2004.
- [67] Han Liu, Larry Wasserman, and John D Lafferty. Exponential concentration for mutual information estimation with application to forests. In *Advances in Neural Information Processing Systems*, pages 2537–2545, 2012.
- [68] Jean-Philippe Lachaux, Eugenio Rodriguez, Jacques Martinerie, and Francisco J Varela. Measuring phase synchrony in brain signals. *Human brain mapping*, 8(4):194–208, 1999.
- [69] Francisco Varela, Jean-Philippe Lachaux, Eugenio Rodriguez, and Jacques Martinerie. The brainweb: phase synchronization and large-scale integration. *Nature reviews neuroscience*, 2(4):229, 2001.
- [70] Michael I Posner. Orienting of attention. *Quarterly journal of experimental psychology*, 32(1):3–25, 1980.
- [71] J Schalkwijk and Thomas Kailath. A coding scheme for additive noise channels with feedback–I: No bandwidth constraint. *Information Theory, IEEE Transactions on*, 12(2):172–182, 1966.
- [72] Amiram Grinvald and Rina Hildesheim. VSDI: a new era in functional imaging of cortical dynamics. *Nature Reviews Neuroscience*, 5(11):874, 2004.

- [73] Christopher J. Watson, B. Jill Venton, and Robert T. Kennedy. In vivo measurements of neurotransmitters by microdialysis sampling. *Analytical Chemistry*, 78(5):1391–1399, Mar 2006.
- [74] Danielle Smith Bassett and Ed Bullmore. Small-world brain networks. *The Neuroscientist*, 12(6):512–523, 2006.
- [75] Sophie Achard, Raymond Salvador, Brandon Whitcer, John Suckling, and ED Bullmore. A resilient, low-frequency, small-world human brain functional network with highly connected association cortical hubs. *Journal of Neuroscience*, 26(1):63–72, 2006.
- [76] Scott J Russo and Eric J Nestler. The brain reward circuitry in mood disorders. *Nature Reviews Neuroscience*, 14(9):609, 2013.
- [77] Juliet Popper Shaffer. Multiple hypothesis testing. *Annual review of psychology*, 46(1):561–584, 1995.
- [78] Boyan Duan, Aaditya Ramdas, Sivaraman Balakrishnan, and Larry Wasserman. Interactive martingale tests for the global null. *arXiv:1909.07339 [stat.ME]*, 2019.
- [79] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1):289–300, 1995.
- [80] Daphne Koller, Nir Friedman, and Francis Bach. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- [81] Nihat Ay and Daniel Polani. Information flows in causal networks. *Advances in Complex Systems*, 11(01):17–41, 2008.
- [82] Dominik Janzing, David Balduzzi, Moritz Grosse-Wentrup, and Bernhard Schölkopf. Quantifying causal influences. *Ann. Statist.*, 41(5):2324–2358, 10 2013.

A Proof of Proposition 1

Proof of Proposition 1. (\Rightarrow) Suppose there exists some $E'_t \in \mathcal{E}'_t$ that has M -information flow. That is,

$$\exists \mathcal{E}''_t \subseteq \mathcal{E}_t \setminus \{E'_t\} \quad \text{s.t.} \quad I(M; X(E'_t) | X(\mathcal{E}''_t)) > 0. \quad (77)$$

Then,

$$I(M; X(\mathcal{E}'_t) | X(\mathcal{E}''_t)) = I(M; X(E'_t) | X(\mathcal{E}''_t)) + I(M; X(\mathcal{E}'_t \setminus \{E'_t\}) | X(\mathcal{E}''_t), X(E'_t)) \quad (78)$$

$$\stackrel{(a)}{\geq} I(M; X(E'_t) | X(\mathcal{E}''_t)) \stackrel{(b)}{>} 0 \quad (79)$$

where (a) follows from the non-negativity of conditional mutual information and (b) from (77). Taking $\mathcal{R}'_t := \mathcal{E}''_t$ in Definition 5, we see that set \mathcal{E}'_t has M -information flow.

(\Leftarrow) Next, suppose that the set \mathcal{E}'_t has M -information flow, as per Definition 5. That is, there exists a set $\mathcal{R}'_t \subseteq \mathcal{E}_t$ such that

$$I(M; X(\mathcal{E}'_t) | X(\mathcal{R}'_t)) > 0. \quad (80)$$

Also, let $\{E_t^{(1)}, E_t^{(2)}, \dots, E_t^{(K)}\}$ be any ordering of the nodes in \mathcal{E}'_t (where $K = |\mathcal{E}'_t|$). Then by the chain rule of mutual information,

$$0 < I(M; X(\mathcal{E}'_t) | X(\mathcal{R}'_t)) \quad (81)$$

$$= \sum_{k=1}^K I\left(M; X(E_t^{(k)}) \mid X(\mathcal{R}'_t), X\left(\bigcup_{j=1}^{k-1} \{E_t^{(j)}\}\right)\right). \quad (82)$$

By the non-negativity of conditional mutual information, at least one of the terms in the summation must be strictly positive. Let the index of this term be k^* . Hence, there exists $E'_t := E_t^{(k^*)}$ and $\mathcal{E}'_t := \mathcal{R}'_t \cup \{E_t^{(1)}, \dots, E_t^{(k^*-1)}\}$, such that

$$I(M; X(E'_t) | X(\mathcal{E}'_t)) > 0. \quad (83)$$

In other words, there exists an edge $E'_t \in \mathcal{E}'_t$ that has M -information flow. \square

B Proof of Proposition 9

Proof of Proposition 9. Consider the set of all $E_t \in \mathcal{E}_t$ that have M -information flow. That is, E_t must satisfy

$$\exists \mathcal{E}'_t \subseteq \mathcal{E}_t \quad \text{s.t.} \quad I(M; X(E_t) | X(\mathcal{E}'_t)) > 0. \quad (84)$$

Define

$$\begin{aligned} \mathcal{R}_t &:= \{E_t \in \mathcal{E}_t : (84) \text{ holds}\}, \\ \mathcal{S}_t &:= \mathcal{E}_t \setminus \mathcal{R}_t. \end{aligned} \quad (85)$$

Then, we claim that \mathcal{R}_t and \mathcal{S}_t satisfy equations (63) and (64).

First, note that if $\mathcal{S}_t \neq \emptyset$, then for every $S_t \in \mathcal{S}_t$, we must have that

$$\forall \mathcal{E}'_t \subseteq \mathcal{E}_t, \quad I(M; X(S_t) | X(\mathcal{E}'_t)) = 0. \quad (86)$$

If not, then $S_t \in \mathcal{R}_t$ by (85), which implies that $S_t \notin \mathcal{S}_t$, which is a contradiction. Hence, we see that no edge in \mathcal{S}_t has M -information flow. Therefore, by Proposition 1, the set \mathcal{S}_t has no M -information flow. This directly implies the condition in (64).

Next, we claim that if $\mathcal{R}_t \neq \emptyset$, then for every $R_t \in \mathcal{R}_t$, if $\mathcal{E}'_t \subseteq \mathcal{E}_t$ is a set that satisfies

$$I(M; X(R_t) | X(\mathcal{E}'_t)) > 0, \quad (87)$$

then $\mathcal{R}'_t := \mathcal{E}'_t \cap \mathcal{R}_t$ satisfies

$$I(M; X(R_t) | X(\mathcal{R}'_t)) > 0. \quad (88)$$

Let $\mathcal{S}'_t := \mathcal{E}'_t \setminus \mathcal{R}'_t$, so that $\mathcal{S}'_t \subseteq \mathcal{S}_t$. Then,

$$I(M; X(R_t) | X(\mathcal{R}'_t), X(\mathcal{S}'_t)) > 0 \quad (89)$$

by (87). So,

$$I(M; X(R_t) | X(\mathcal{R}'_t)) \stackrel{(a)}{=} I(M; X(R_t), X(\mathcal{S}'_t) | X(\mathcal{R}'_t)) - I(M; X(\mathcal{S}'_t) | X(\mathcal{R}'_t), X(R_t)) \quad (90)$$

$$\stackrel{(b)}{=} I(M; X(R_t), X(\mathcal{S}'_t) | X(\mathcal{R}'_t)) \quad (91)$$

$$\stackrel{(c)}{=} I(M; X(R_t) | X(\mathcal{R}'_t), X(\mathcal{S}'_t)) + I(M; X(\mathcal{S}'_t) | X(\mathcal{R}'_t)) \quad (92)$$

$$\stackrel{(d)}{=} I(M; X(R_t) | X(\mathcal{R}'_t), X(\mathcal{S}'_t)) \quad (93)$$

$$\stackrel{(e)}{>} 0,$$

where (a) and (c) follow from the chain rule, (b) and (d) follow from (64), and (e) follows from (89). Thus, condition (63) also holds. \square

C Synergistic Information Flow

C.1 Partial Information Decomposition preliminaries

The literature on Partial Information Decomposition seeks to find a decomposition for the mutual information between a message, M , and a set of random variables, $\{X_1, X_2, \dots\}$ into several individually meaningful, non-negative terms [50]. For our purposes, it suffices to consider the *bivariate* case, i.e., the decomposition of $I(M; X, Y)$ into non-negative components. In the bivariate case, it is well-understood *how many* components there ought to be, and what these quantities *intuitively represent*, but as yet, there is no consensus on a single set of definitions [50].

There is, however, consensus on a basic set of properties that we expect these components to satisfy. For our purposes, we will only make use of the basic properties stated here, so that *any definition* of the aforementioned components which satisfies these properties suffices for our theory.

In the bivariate case, the mutual information between M and (X, Y) is decomposed into four components: information about M which is (i) unique to X and not present in Y , (ii) unique to Y and not present in X , (iii) redundantly present in both X and Y , and (iv) synergistically present in X and Y . In the notation of [49], the decomposition is written as:

$$I(M; (X, Y)) = UI(M : X \setminus Y) + UI(M : Y \setminus X) + SI(M : X; Y) + CI(M : X; Y), \quad (94)$$

where the components are ordered exactly as stated above. Note that SI refers to “shared”, and hence redundant, information, while CI refers to “complementary”, and hence synergistic, information. We shall continue to use the terms “redundant” and “synergistic”, however, since they are more meaningful in this context. Also, in what follows, we shall assume that SI and CI are symmetric in X and Y . This is usually an additional condition that is imposed when defining these quantities, but here, we take it as given.

Given what we want the four components to represent, we would also expect the following to hold:

$$\begin{aligned} I(M; X) &= UI(M : X \setminus Y) + SI(M : X; Y), \\ I(M; Y) &= UI(M : Y \setminus X) + SI(M : X; Y). \end{aligned} \quad (95)$$

As a natural consequence, this means that the conditional mutual information will satisfy:

$$\begin{aligned} I(M; X | Y) &= I(M; (X, Y)) - I(M; Y) \\ &= UI(M : X \setminus Y) + CI(M : X; Y), \\ I(M; Y | X) &= I(M; (Y, X)) - I(M; X) \\ &= UI(M : Y \setminus X) + CI(M : X; Y). \end{aligned} \quad (96)$$

Finally, we want each of these components to always be non-negative:

$$\begin{aligned} UI(M : X \setminus Y) &\geq 0 & SI(M : X; Y) &\geq 0 \\ UI(M : Y \setminus X) &\geq 0 & CI(M : X; Y) &\geq 0. \end{aligned} \quad (97)$$

It is not obvious that a consistent definition of these four quantities which also satisfies the equations stated above even *exists*, but in fact, additional properties are required to obtain a unique definition. For instance, see [49] for one such development.

As stated before, our theory only relies on the properties stated in this section. As a result, our theorem on the equivalence of information flow definitions holds irrespective of what definition is used, exactly, for synergistic information. It only matters that the definition used satisfies the basic properties presented here.

C.2 Equivalence of information flow definitions

Proof of Proposition 2. (\Rightarrow) Suppose the edge E_t has strictly positive M -information flow. Then,

$$\exists \mathcal{E}'_t \subseteq \mathcal{E}_t \quad \text{s.t.} \quad I(M; X(E_t) | X(\mathcal{E}'_t)) > 0. \quad (98)$$

If $I(M; X(E_t)) > 0$ with $\mathcal{E}'_t = \emptyset$ in (98), then condition 1 in Definition 6 holds, so nothing remains to be shown. If not, then $I(M; X(E_t)) = 0$, so (98) implies that there must exist some $\mathcal{E}'_t \neq \emptyset$ such that

$$I(M; X(E_t) | X(\mathcal{E}'_t)) > 0, \quad (99)$$

which, by (96), is equivalent to

$$UI(M : X(E_t) \setminus X(\mathcal{E}'_t)) + CI(M : X(E_t); X(\mathcal{E}'_t)) > 0. \quad (100)$$

However, since $I(M; X(E_t)) = 0$, we must have $UI(M : X(E_t) \setminus X(\mathcal{E}'_t)) = 0$ by (95) and (97). Hence,

$$\exists \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\} \quad \text{s.t.} \quad CI(M : X(E_t); X(\mathcal{E}'_t)) > 0. \quad (101)$$

So the implication in the forward direction holds.

(\Leftarrow) For the converse, suppose that E_t has no M -information flow. That is,

$$I(M; X(E_t) | X(\mathcal{E}'_t)) = 0 \quad \forall \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\}. \quad (102)$$

By (96), this implies that

$$UI(M : X(E_t) \setminus X(\mathcal{E}'_t)) + CI(M : X(E_t); X(\mathcal{E}'_t)) = 0 \quad \forall \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\}. \quad (103)$$

Since UI and CI are both non-negative by (97), we must have that

$$CI(M : X(E_t); X(\mathcal{E}'_t)) = 0 \quad \forall \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\}. \quad (104)$$

This proves the converse. \square

D Miscellaneous Proofs from Section 5

D.1 Proof of Lemma 10

Proof of Lemma 10. Consider a subset of hidden nodes $\mathcal{H}'_t \subseteq \mathcal{H}_t$ that is not M -relevant. Then, by Definition 14, $\mathbb{Q}(\mathcal{H}'_t)$ carries no M -information flow in \mathcal{G} . This means that

$$\forall \mathcal{E}'_t \subseteq \mathcal{E}_t, \quad I(M; X(\mathbb{Q}(\mathcal{H}'_t)) | X(\mathcal{E}'_t)) = 0. \quad (105)$$

Specifically, taking $\mathcal{E}'_t = \tilde{\mathcal{E}}_t$, we have

$$I(M; X(\mathbb{Q}(\mathcal{H}'_t)) | X(\tilde{\mathcal{E}}_t)) = 0. \quad (106)$$

Therefore, by Definition 15, \mathcal{H}'_t is M -derived. Thus, if \mathcal{H}'_t is *not* M -relevant, it *is* M -derived. Taking the contrapositive, if \mathcal{H}_t is *not* M -derived, then it *is* M -relevant. \square

D.2 Proof of proposition 11

Proof of Proposition 11. We are given that

$$I(M; X(\tilde{\mathcal{E}}_{t+1}) \mid X(\tilde{\mathcal{E}}_t)) > 0, \quad (107)$$

and must prove that the hidden nodes at time t , \mathcal{H}_t , are *not* M -derived.

First note that, since $\mathbb{Q}(\tilde{\mathcal{V}}_{t+1}) = \tilde{\mathcal{E}}_{t+1} \cup (\tilde{\mathcal{V}}_{t+1} \times \mathcal{H}_{t+2})$, we must have

$$I(M; X(\mathbb{Q}(\tilde{\mathcal{V}}_{t+1})) \mid X(\tilde{\mathcal{E}}_t)) = I(M; X(\tilde{\mathcal{E}}_{t+1}), X(\tilde{\mathcal{V}}_{t+1} \times \mathcal{H}_{t+2}) \mid X(\tilde{\mathcal{E}}_t)) \quad (108)$$

$$= I(M; X(\tilde{\mathcal{E}}_{t+1}) \mid X(\tilde{\mathcal{E}}_t)) + I(M; X(\tilde{\mathcal{V}}_{t+1} \times \mathcal{H}_{t+2}) \mid X(\tilde{\mathcal{E}}_{t+1}), X(\tilde{\mathcal{E}}_t)) \quad (109)$$

$$\geq I(M; X(\tilde{\mathcal{E}}_{t+1}) \mid X(\tilde{\mathcal{E}}_t)) \quad (110)$$

$$> 0, \quad (111)$$

where the last line follows from the fact that conditional mutual information is non-negative, and from (107).

Next, observe that Local Markovity conditions (Proposition 6) *must* hold on the *entire* graph \mathcal{G} , which consists of both observed and hidden nodes. If we apply the Local Markovity condition to $\tilde{\mathcal{V}}_{t+1}$, we have $M - X(\mathcal{P}(\tilde{\mathcal{V}}_{t+1})) - X(\mathbb{Q}(\tilde{\mathcal{V}}_{t+1}))$, or in other words

$$I(M; X(\mathbb{Q}(\tilde{\mathcal{V}}_{t+1})) \mid X(\mathcal{P}(\tilde{\mathcal{V}}_{t+1}))) = 0. \quad (112)$$

Note that $\mathcal{P}(\tilde{\mathcal{V}}_{t+1}) = \tilde{\mathcal{E}}_t \cup \tilde{\mathbb{Q}}(\mathcal{H}_t)$, where $\tilde{\mathbb{Q}}(\mathcal{H}_t) := \mathcal{H}_t \times \tilde{\mathcal{V}}_{t+1}$ is the subset comprising outgoing edges of \mathcal{H}_t that go to $\tilde{\mathcal{V}}_{t+1}$. Therefore,

$$I(M; X(\mathbb{Q}(\tilde{\mathcal{V}}_{t+1})) \mid X(\tilde{\mathcal{E}}_t), X(\tilde{\mathbb{Q}}(\mathcal{H}_t))) = 0. \quad (113)$$

Expanding this conditional mutual information, we get

$$I(M; X(\mathbb{Q}(\tilde{\mathcal{V}}_{t+1}), X(\tilde{\mathbb{Q}}(\mathcal{H}_t)) \mid X(\tilde{\mathcal{E}}_t)) - I(M; X(\tilde{\mathbb{Q}}(\mathcal{H}_t)) \mid X(\tilde{\mathcal{E}}_t)) = 0. \quad (114)$$

So we have

$$I(M; X(\tilde{\mathbb{Q}}(\mathcal{H}_t)) \mid X(\tilde{\mathcal{E}}_t)) = I(M; X(\mathbb{Q}(\tilde{\mathcal{V}}_{t+1}), X(\tilde{\mathbb{Q}}(\mathcal{H}_t)) \mid X(\tilde{\mathcal{E}}_t)) \quad (115)$$

$$= I(M; X(\mathbb{Q}(\tilde{\mathcal{V}}_{t+1})) \mid X(\tilde{\mathcal{E}}_t)) + I(M; X(\tilde{\mathbb{Q}}(\mathcal{H}_t)) \mid X(\mathbb{Q}(\tilde{\mathcal{V}}_{t+1}), X(\tilde{\mathcal{E}}_t)) > 0, \quad (116)$$

where the final inequality follows from (111) and the fact that conditional mutual information is non-negative. Finally, since $\tilde{\mathbb{Q}}(\mathcal{H}_t) \subset \mathbb{Q}(\mathcal{H}_t)$, we have that $I(M; X(\mathbb{Q}(\mathcal{H}_t)) \mid X(\tilde{\mathcal{E}}_t)) > 0$, just as we in equations (108)–(111). Hence, the Markov chain $M - X(\tilde{\mathcal{E}}_t) - X(\mathbb{Q}(\mathcal{H}_t))$ does not hold, so by Definition 15, \mathcal{H}_t are not M -derived. \square

E On the Uniqueness of Our Definition of Information Flow

From the perspective of designing an axiomatic framework, it is desirable to find a minimal set of properties that gives rise to a unique definition of information flow. Although Property 1 helped us motivate a definition for information flow, it did not uniquely specify a definition. Indeed, the all-zero function as well as the all-ones function also satisfy the property, although they are not particularly useful definitions of information flow.

In this section, we provide a set of properties that uniquely leads to our definition of information flow. However, we must acknowledge that we arrived at these properties with the benefit of hindsight, after having proved many other properties of our definition. As such, they are mathematically very similar to

our definition, and one might feel uncomfortable with the idea of imposing such a set of properties at the very outset. Our goal here is only to begin a discussion in this direction: a search for a more abstract set of properties that leads to a unique definition of information flow would be a worthy endeavour in future.

Property 4. *Let \mathcal{E} be a computational system, and let $\mathcal{F}_M : \mathcal{E} \rightarrow \{0, 1\}$ be an indicator of the presence of information flow about M on an edge. That is, $\mathcal{F}_M(E) = 1$, if information about M flows on the edge $E \in \mathcal{E}$, and $\mathcal{F}_M(E) = 0$ otherwise. We now state three conditions \mathcal{F}_M must satisfy, which naturally leads to our definition of information flow (Definition 4):*

$$4a) I(M; X(E_t)) > 0 \quad \Rightarrow \quad \mathcal{F}_M(E_t) = 1,$$

$$4b) \exists \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\} \quad \text{s.t.} \quad I(M; X(\mathcal{E}'_t) | X(E_t)) > I(M; X(\mathcal{E}'_t)) \quad \Rightarrow \quad \mathcal{F}_M(E_t) = 1,$$

$$4c) I(M; X(E_t) | X(\mathcal{E}'_t)) = 0 \quad \forall \mathcal{E}'_t \subseteq \mathcal{E}_t \quad \Rightarrow \quad \mathcal{F}_M(E_t) = 0.$$

Property 4a is a very natural and intuitive requirement for information flow. Property 4b states that an edge should be considered to carry information about M , if upon conditioning, its transmission *increases* the information that some set $X(\mathcal{E}'_t)$ conveys about M . Property 4c is reminiscent of the separability property from Proposition 9, and states that if an edge has no dependence with M , no matter what other transmission is conditioned upon, then it can carry no information flow about M .

Effectively, Property 4a states that if an edge has unique or redundant information about M , then it must carry information flow, while Property 4b states that if an edge has synergistic information about M along with some other set of transmissions, then it must carry information flow. Finally, Property 4c states that if all three of these components are absent, then that edge carries no information flow. This also explains how, if any one of these three properties is absent, our definition is no longer unique.

As we acknowledged previously, some of these properties could be seen as too restrictive or contrived, and a more abstract set of properties is certainly desirable. Nevertheless, these properties do uniquely identify our definition of information flow.

Proposition 12 (Uniqueness). *If \mathcal{F}_M is an indicator of information flow that satisfies the conditions in Property 4, then $\mathcal{F}_M(E_t) = 1$ if and only if E_t has M -information flow, per Definition 4.*

Proof. (\Rightarrow) Suppose the edge E_t has no M -information flow per Definition 4. This directly implies the condition in Property 4c. Hence, $\mathcal{F}_M(E_t) = 0$. This proves that if $\mathcal{F}_M(E_t) = 1$, the edge E_t must have M -information flow.

(\Leftarrow) Suppose the edge E_t has M -information flow per Definition 4. Then,

$$\exists \mathcal{E}'_t \subseteq \mathcal{E}_t \setminus \{E_t\} \quad \text{s.t.} \quad I(M; X(E_t) | X(\mathcal{E}'_t)) > 0. \quad (117)$$

If $\mathcal{E}'_t = \emptyset$, $I(M; X(E_t)) > 0$, so by Property 4a, $\mathcal{F}_M(E_t) = 1$. If $I(M; X(E_t)) = 0$, then (117) guarantees the existence of some $\mathcal{E}'_t \neq \emptyset$ such that

$$I(M; X(E_t) | X(\mathcal{E}'_t)) > 0 \quad (118)$$

$$\Rightarrow I(M; X(\mathcal{E}'_t)) + I(M; X(E_t) | X(\mathcal{E}'_t)) \stackrel{(a)}{>} I(M; X(\mathcal{E}'_t)) \quad (119)$$

$$\Rightarrow I(M; X(E_t), X(\mathcal{E}'_t)) \stackrel{(b)}{>} I(M; X(\mathcal{E}'_t)) \quad (120)$$

$$\Rightarrow I(M; X(E_t)) + I(M; X(\mathcal{E}'_t) | X(E_t)) \stackrel{(c)}{>} I(M; X(\mathcal{E}'_t)) \quad (121)$$

$$\Rightarrow I(M; X(\mathcal{E}'_t) | X(E_t)) \stackrel{(d)}{>} I(M; X(\mathcal{E}'_t)), \quad (122)$$

where in (a), we simply added $I(M; X(\mathcal{E}'_t))$ to both sides; in (b) and (c), we used the chain rule in two different ways; and in (d), we used the fact that $I(M; X(E_t)) = 0$. So, by Property 4b, we have that $\mathcal{F}_M(E_t) = 1$. This proves the converse. \square

Remark It should be noted that Definition 4 only specifies *whether or not* a given edge has M -information flow. It does not *quantify* this flow. So Proposition 12 demonstrates the uniqueness of our definition up to an unspecified information volume. If we require that the conditions in Property 4 hold, then any quantitative definition of information flow will go to zero at an edge if and only if the M -information flow carried by that edge is zero.

F Miscellaneous Derivations from Section 6

F.1 Derivation of Expressions in the Second FFT Example from Section 6.2

Here, we derive the expressions used in Figure 13. Recall that $Y_i = \omega^{-iM}/4$, where $\omega = e^{-j2\pi/4} = -j$.

$$Y_{02} = Y_0 + Y_2 = \frac{1}{4} + \frac{\omega^{-2M}}{4} = \frac{1}{4}(1 + \omega^{-2M}) \quad (123)$$

$$Y_{13} = Y_1 + Y_3 = \frac{\omega^{-M}}{4} + \frac{\omega^{-3M}}{4} = \frac{\omega^{-M}}{4}(1 + \omega^{-2M}) \quad (124)$$

$$Y'_{02} = Y_0 + \omega^2 Y_2 = \frac{1}{4} + (-1)\frac{\omega^{-2M}}{4} = \frac{1}{4}(1 - \omega^{-2M}) \quad (125)$$

$$Y'_{13} = Y_1 + \omega^2 Y_3 = \frac{\omega^{-M}}{4} + (-1)\frac{\omega^{-3M}}{4} = \frac{\omega^{-M}}{4}(1 - \omega^{-2M}) \quad (126)$$

Next, we show that these intermediate values actually yield the expected values of \tilde{Y} .

$$\tilde{Y}_0 = Y_{02} + Y_{13} = \frac{1}{4}(1 + \omega^{-2M} + \omega^{-M} + \omega^{-3M}) \quad (127)$$

$$= \begin{cases} \frac{1}{4}(1 + 1 + 1 + 1), & M = 0 \\ \frac{1}{4}(1 + j + j^2 + j^3), & M = 1 \end{cases} \quad (128)$$

$$= 1 - M \quad (129)$$

$$\tilde{Y}_1 = Y'_{02} + \omega Y'_{13} = \frac{1}{4}(1 - \omega^{-2M} + \omega^{1-M} - \omega^{1-3M}) \quad (130)$$

$$= \begin{cases} \frac{1}{4}(1 - 1 + \omega - \omega), & M = 0 \\ \frac{1}{4}(1 - j^2 + 1 - j^2), & M = 1 \end{cases} \quad (131)$$

$$= M \quad (132)$$

$$\tilde{Y}_2 = Y_{02} + \omega^2 Y_{13} = \frac{1}{4}(1 + \omega^{-2M} + \omega^{2-M} + \omega^{2-3M}) \quad (133)$$

$$= \frac{1}{4}(1 + \omega^{-2M} - \omega^{-M} - \omega^{-3M}) \quad (134)$$

$$= \begin{cases} \frac{1}{4}(1 + 1 - 1 - 1), & M = 0 \\ \frac{1}{4}(1 - 1 - \omega^{-1} + \omega^{-1}), & M = 1 \end{cases} \quad (135)$$

$$= 0 \quad (136)$$

$$\tilde{Y}_3 = Y'_{02} + \omega^3 Y'_{13} = \frac{1}{4}(1 - \omega^{-2M} + \omega^{3-M} - \omega^{3-3M}) \quad (137)$$

$$= \frac{1}{4}(1 - \omega^{-2M} + \omega^3(\omega^{-M} - \omega^{-3M})) \quad (138)$$

$$= \begin{cases} \frac{1}{4}(1 - 1 - \omega(1 - 1)), & M = 0 \\ \frac{1}{4}(1 - (-1) - \omega(\omega^{-1} + \omega^{-1})), & M = 1 \end{cases} \quad (139)$$

$$= 0 \quad (140)$$

F.2 Derivation of the Markov Chain Failure in Section 6.3

We wish to show that in the canonical example from Section 6.3, $M - [\widehat{M}_1, \widehat{M}_2] - (M - \widehat{M}_2)$ is not a valid Markov chain. Recall that $Z_1, Z_2, Z_3 \sim \text{i.i.d. } \mathcal{N}(0, \sigma^2)$ and $M \sim \mathcal{N}(0, 1)$. Let $h(\cdot)$ denote differential entropy. Then,

$$I(M; M - \widehat{M}_2, \widehat{M}_2) > I(M; \widehat{M}_3) = h(\widehat{M}_3) - h(\widehat{M}_3 | M) \quad (141)$$

$$= \frac{1}{2} \log \left(2\pi e \left(1 + \frac{\sigma^2}{3} \right) \right) - \frac{1}{2} \log \left(2\pi e \left(\frac{\sigma^2}{3} \right) \right) \quad (142)$$

$$= \frac{1}{2} \log \left(1 + \frac{3}{\sigma^2} \right). \quad (143)$$

Here, we have used the fact that if $Y \sim \mathcal{N}(0, \sigma^2)$ is a zero-mean multivariate Gaussian random variable with variance σ^2 , then its differential entropy is given by [55, Thm. 8.4.1]

$$h(Y) = \frac{1}{2} \log(2\pi e \sigma^2) \text{ nats.} \quad (144)$$

Next, note that since $\widehat{M}_1 = M + Z_1$ and $\widehat{M}_2 = M + \frac{1}{2}(Z_1 + Z_2)$, \widehat{M}_1 has no extra information about M , given \widehat{M}_2 . This is obvious when we think of \widehat{M}_1 as being $\widehat{M}_1 = \widehat{M}_2 + Z'$, where $Z' = \frac{1}{2}(Z_1 - Z_2)$, and it can be shown that $Z' \perp \widehat{M}_2$:

$$\mathbb{E}[\widehat{M}_2 Z'] = \mathbb{E} \left[\left(M + \frac{1}{2}(Z_1 + Z_2) \right) Z' \right] \quad (145)$$

$$= \mathbb{E}[M Z'] + \frac{1}{4} \mathbb{E}[(Z_1 + Z_2)(Z_1 - Z_2)] \quad (146)$$

$$= 0 + \frac{1}{4} \mathbb{E}[Z_1^2 - Z_2^2] \quad (147)$$

$$= \frac{1}{4}(\sigma^2 - \sigma^2) = 0. \quad (148)$$

Since all variables involved are zero-mean Gaussians, this naturally implies that $\widehat{M}_2 \perp Z'$. Thus, from our previous argument, \widehat{M}_1 has no extra information about M when given \widehat{M}_2 , or in other words, $M - \widehat{M}_2 - \widehat{M}_1$ is a valid Markov chain. Therefore,

$$I(M; \widehat{M}_1, \widehat{M}_2) = I(M; \widehat{M}_2) + I(M; \widehat{M}_1 | \widehat{M}_2) \quad (149)$$

$$= \frac{1}{2} \log \left(1 + \frac{2}{\sigma^2} \right) + 0, \quad (150)$$

derived in the same way as (143). From (143) and (150), we can conclude that $I(M; \widehat{M}_3) > I(M; \widehat{M}_2)$, and therefore

$$I(M; M - \widehat{M}_2, \widehat{M}_2) > I(M; \widehat{M}_1, \widehat{M}_2) \quad (151)$$

$$I(M; M - \widehat{M}_2, \widehat{M}_2, \widehat{M}_1) > I(M; \widehat{M}_1, \widehat{M}_2) \quad (152)$$

$$I(M; M - \widehat{M}_2, \widehat{M}_2, \widehat{M}_1) - I(M; \widehat{M}_1, \widehat{M}_2) > 0 \quad (153)$$

$$I(M; M - \widehat{M}_2 | \widehat{M}_1, \widehat{M}_2) > 0. \quad (154)$$

Thus, the stated Markov chain, $M - [\widehat{M}_1, \widehat{M}_2] - (M - \widehat{M}_2)$, cannot hold.