

Directed Data-Processing Inequalities for Systems with Feedback

Derpich, Milan; Østergaard, Jan

Published in:
Entropy

DOI (link to publication from Publisher):
[10.3390/e23050533](https://doi.org/10.3390/e23050533)

Creative Commons License
CC BY 4.0

Publication date:
2021

Document Version
Publisher's PDF, also known as Version of record

[Link to publication from Aalborg University](#)

Citation for published version (APA):
Derpich, M., & Østergaard, J. (2021). Directed Data-Processing Inequalities for Systems with Feedback. *Entropy*, 23(5), Article 533. <https://doi.org/10.3390/e23050533>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal -

Take down policy

If you believe that this document breaches copyright please contact us at vbn@aub.aau.dk providing details, and we will remove access to the work immediately and investigate your claim.

Article

Directed Data-Processing Inequalities for Systems with Feedback

Milan S. Derpich ^{1,*}  and Jan Østergaard ^{2,*} 
¹ Department of Electronic Engineering, Universidad Técnica Federico Santa María, Av. España 1680, Valparaíso 2390123, Chile

² Department of Electronic Systems, Aalborg University, 9220 Aalborg, Denmark

* Correspondence: milan.derpich@usm.cl (M.S.D.); jo@es.aau.dk (J.Ø.)

Abstract: We present novel data-processing inequalities relating the mutual information and the directed information in systems with feedback. The internal deterministic blocks within such systems are restricted only to be causal mappings, but are allowed to be non-linear and time varying, and randomized by their own external random input, can yield any stochastic mapping. These randomized blocks can for example represent source encoders, decoders, or even communication channels. Moreover, the involved signals can be arbitrarily distributed. Our first main result relates mutual and directed information and can be interpreted as a law of conservation of information flow. Our second main result is a pair of data-processing inequalities (one the conditional version of the other) between nested pairs of random sequences entirely within the closed loop. Our third main result introduces and characterizes the notion of in-the-loop (ITL) transmission rate for channel coding scenarios in which the messages are internal to the loop. Interestingly, in this case the conventional notions of transmission rate associated with the entropy of the messages and of channel capacity based on maximizing the mutual information between the messages and the output turn out to be inadequate. Instead, as we show, the ITL transmission rate is the unique notion of rate for which a channel code attains zero error probability if and only if such an ITL rate does not exceed the corresponding directed information rate from messages to decoded messages. We apply our data-processing inequalities to show that the supremum of achievable (in the usual channel coding sense) ITL transmission rates is upper bounded by the supremum of the directed information rate across the communication channel. Moreover, we present an example in which this upper bound is attained. Finally, we further illustrate the applicability of our results by discussing how they make possible the generalization of two fundamental inequalities known in networked control literature.

Keywords: data-processing inequality; directed information; mutual information; networked control; feedback capacity



Citation: Derpich, M.S.; Østergaard, J. Directed Data-Processing Inequalities for Systems with Feedback. *Entropy* **2021**, *23*, 533. <https://doi.org/10.3390/e23050533>

Academic Editor: José A. Tenreiro Machado

Received: 24 March 2021

Accepted: 21 April 2021

Published: 26 April 2021

Publisher's Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

1. Introduction

The data-processing inequality states that if x, y, z are random variables such that x and z become independent when conditioning upon y , then:

$$I(x; y) \geq I(x; z) \quad (1)$$

$$I(y; z) \geq I(x; z), \quad (2)$$

where $I(x; y)$ denotes the mutual information between x and y [1] (p. 252) (a definition of mutual information is provided in Section 2.2 below). Among its many uses, the data-processing inequality plays a key role in the proof of the converse part (i.e., outer bounds) in rate-distortion [1–5] (p. 317), channel capacity [1,6–8] (pp. 208, 217, 540 and 566), and joint source-channel coding theorems [1,9–11] (p. 221), and has recently been extended to von Neumann algebras, which have applications in quantum field theory (see [12] and the references therein).

It is well known that mutual information has an important limitation in systems with feedback, such as the one shown in Figure 1a. In this system, p, q, r, s, e, u, x , and y are random sequences, and the blocks S_1, \dots, S_4 are deterministic causal mappings with an added delay of at least one sample. These blocks, randomized by their exogenous random inputs p, q, r, s , may yield any causal stochastic dynamic mappings. As pointed out in [13], for sequences inside the loop, such as x and y , $I(x; y)$ does not distinguish the probabilistic interdependence produced by the effect x has on y from that stemming from the influence of y on x . This limitation motivated the introduction of the directed information in [13]. This notion assesses the amount of information that causally “flows” from a given random and ordered sequence to another. For this reason, it has increasingly found use in diverse applications, including characterizing the capacity of channels with feedback [13–16], the rate distortion function under causality constraints [5], establishing some of the fundamental limitations in networked control [17–23], determining causal relationships in neural networks [24], and portfolio theory and hypothesis testing [25], to name a few.

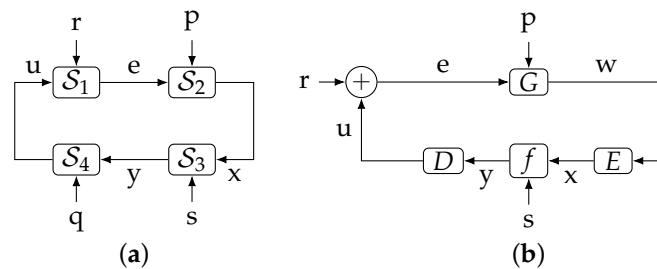


Figure 1. (a): The general system considered in this work. (b): A special case of (a), corresponding to the closed-loop system studied in [18].

The directed information from a random sequence x^k to a random sequence y^k is defined as:

$$I(x^k \rightarrow y^k) \triangleq \sum_{i=1}^k I(y(i); x^i | y^{i-1}), \quad (3)$$

where the notation x^i represents the sequence $x(1), x(2), \dots, x(i)$ and $I(x; y | z)$ is the mutual information between x and y conditioned on (or given) z (hereafter we use non-italic letters, such as x , for random variables, denoting a particular realization by the corresponding italic character, x). The causality inherent in this definition becomes evident when comparing it with the mutual information between x^k and y^k , given by $I(x^k; y^k) = \sum_{i=1}^k I(y(i); x^k | y^{i-1})$. In the latter sum, what matters is the amount of information about the *entire* sequence x^k present in $y(i)$, given the past values y^{i-1} . By contrast, in the conditional mutual information in the sum of (3), only the past and current values of x^k are considered, that is, x^i . Thus, $I(x^k \rightarrow y^k)$ represents the amount of information causally conveyed from x^k to y^k . A related notion is the causally conditioned directed information introduced in [14], defined as:

$$I(x^k \rightarrow y^k \parallel q^k) \triangleq \sum_{i=1}^k I(y(i); x^i | y^{i-1}, q^i). \quad (4)$$

In this paper, we derive inequalities involving directed and mutual information within feedback systems. For this purpose, we consider the general feedback system shown in Figure 1a. In this diagram, the blocks S_1, \dots, S_4 represent possibly non-linear and time-varying causal discrete-time systems such that the total delay of the loop is at least one sample. These blocks can model, for example, source encoders, decoders or even communication channels. In the same figure, r, p, s, q are exogenous random signals (scalars, vectors, or sequences), which could represent, for example, any combination of disturbances, noises,

random initial states, or side information. We note that any of these exogenous signals, in combination with their corresponding deterministic mapping S_i , can also yield any desired stochastic causal mapping (for example, a noisy communication channel, a zero-delay source coder or decoder, or a causal dynamic system with disturbances and a random initial state).

1.1. Main Contributions

1. Our first main contribution is the following theorem. It states a fundamental result, which relates the directed information between two signals within a feedback loop, say x and y , to the mutual information between an external set of signals and y :

Theorem 1. *In the system shown in Figure 1a, it holds that:*

$$I(x^k \rightarrow y^k) = I(q^k, r^k, p^k \rightarrow y^k) - I(q^k, r^k, p^k \rightarrow y^k \parallel x^k) \leq I(p^k, q^k, r^k; y^k), \quad \forall k \in \mathbb{N}, \quad (5)$$

with equality achieved if s is independent of (p, q, r) .

The proof is in Section 3. This fundamental result, which for the cases in which $s \perp (p, q, r)$ can be understood as a *law of conservation of information flow*, is illustrated in Figure 2. (Here, and in the sequel, we use the notation $x \perp y$ to mean “ x is independent of y ”.) For such a case, the information causally conveyed from x to y equals the information flow from (q, r, p) to y . When (p, q, r) are not independent of s , part of the mutual information between (p, q, r) and y (corresponding to the term $I(q^k, r^k, p^k \rightarrow y^k \parallel x^k)$) can be thought of as being “leaked” through s , thus bypassing the forward link from x to y . This provides an intuitive interpretation for (5).

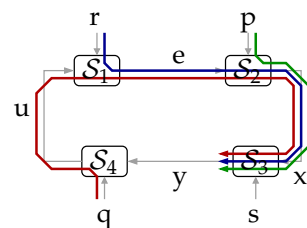


Figure 2. The flow of information between exogenous signals (p, q, r) and the internal signal y equals the directed information from x^k to y^k when $s \perp (p, q, r)$.

Remark 1. *Theorem 1 implies that $I(x^k \rightarrow y^k)$ is only a part of (or at most equal to) the information “flow” between all the exogenous signals entering the loop outside the link $x \rightarrow y$ (namely (q, r, p)), and y . In particular, if (p, q, r) were deterministic, then $I(x^k \rightarrow y^k) = 0$, regardless of the blocks S_1, \dots, S_4 and irrespective of the nature of s .*

2. Our second main result is the following theorem, which relates directed information involving four different sequences internal to the loop. The proof is in Appendix A on page 19.

Theorem 2 (Full Closed-Loop Directed Data-Processing Inequality). *Consider the system shown in Figure 1a.*

- (a) *If $(q, s) \perp (r, p)$ and $q \perp s$, or if $(p, s) \perp (r, q)$ and $p \perp s$, then:*

$$I(x^k \rightarrow y^k) \geq I(e^k \rightarrow u^k). \quad (6)$$

- (b) *If $(q, s) \perp (r, p)$ and $q_{i+1}^k \leftrightarrow q^i \leftrightarrow s^i$ for $i = 1, 2, \dots, k-1$, then :*

$$I(x^k \rightarrow y^k \parallel q^k) \geq I(e^k \rightarrow u^k). \quad (7)$$

(The Markov chain notation $t \leftrightarrow v \leftrightarrow w$ means “ t and w are independent when v is given”). To the best of our knowledge, Theorem 2 is the first result available in the literature providing a lower bound to the gap between two instances of nested directed information, involving four different signals inside the feedback loop. This result can be seen as the first full extension of the open-loop (traditional) data-processing inequality, to arbitrary closed-loop scenarios. (Notice that there is no need to consider systems with more than four mappings, since all external signals entering the loop between a given pair of internal signals can be regarded as exogenous inputs to a single equivalent deterministic mapping.)

3. Our third main contribution is introducing the notion of *in-the-loop* (ITL) transmission rate (in Section 6) for the (seldom considered) channel-coding scenario in which the messages to be transmitted and the communication channel are internal to a feedback loop. We show that the supremum of the directed information rate across such a channel upper bounds the achievable ITL transmission rates. Moreover, we present an example in which this upper bound is attainable. This gives further operational meaning to the directed information rate in closed-loop scenarios.
4. Finally, we provide additional examples of the applicability of our results by discussing how they allow one to obtain the generalizations of two fundamental inequalities known in networked control literature. The first one appears in [18] (Lemma 4.1) and is written in (12) below. This generalization is a consequence of Theorem 4 and is discussed in Remarks 3 and 5 below. The second generalization applies to [20] (Theorem 4.1) and is described on page 6 below. It is an application of Theorem 2 that has just been carried out by the authors in [26], which is all the more important since, as we also reveal in that note, there is a flaw in the proof of [20] (Theorem 4.1).

A key ingredient in proving most of our theorems is provided by Lemma 1, stated in Section 2.5. It allows one to rigorously establish some of the non-trivial conditional independencies that arise in a feedback loop with several (possibly stochastic) dynamic systems.

Put together, the law of conservation of information flow from Theorem 1, our extension of the data processing inequality to general feedback systems from Theorem 2, and our other results constitute both a conceptual framework and a toolbox for addressing information flow problems in feedback systems. We are convinced that this contribution will be instrumental in establishing new results on, e.g., rate-distortion and channel capacity problems with feedback.

The literature review presented next will allow the reader to further assess the novelty and relevance of our results.

1.2. Existing Related Results

There exist several results characterizing the relationship between $I(x^k \rightarrow y^k)$ and $I(x^k; y^k)$. First, it is well known that $I(x^k \rightarrow y^k) \leq I(x^k; y^k)$, with equality if and only if y^k is causally related to x^k [13]. A conservation law of mutual and directed information has been found in [27], which asserts that $I(x^k \rightarrow y^k) + I(0 * y^{k-1} \rightarrow x^k) = I(x^k; y^k)$, where $0 * y^{k-1}$ denotes the concatenation $0, y(1), \dots, y^{k-1}$.

Given its prominence in settings involving feedback, it is perhaps in these scenarios where the directed information becomes most important. For instance, the directed information has been instrumental in characterizing the capacity of channels with feedback (see, e.g., [15,16,28] and the references therein), as well as the rate-distortion function in setups involving feedback [5,20–22,29].

For the simple case in which all the systems $\{\mathcal{S}_i\}_{i=1}^4$ are linear time invariant (LTI) and stable, and assuming $p, x, q = 0$ (deterministically), it was shown in [30] that $I(r^k \rightarrow e^k)$ does not depend on whether there is feedback from e to u or not.

Inequalities between mutual and directed information in a less restricted setup, shown in Figure 1b, have been found in [18,19]. In that setting (a networked-control system), G is a strictly causal LTI dynamic system having (vector) state sequence $\{x(i)\}_{i=0}^\infty$, with $p \triangleq x(0)$ being the random initial state in its state-space representation. The external signal r (which

could correspond to a disturbance) is statistically independent of s , the latter corresponding to, for example, side information or channel noise. Both are also statistically independent of p .

The blocks labeled E , D , and f correspond to an encoder, a decoder, and a channel, respectively, all of which are causal. The channel f maps s^k and x^k to $y(k)$ in a possibly time-varying manner, i.e., $y(k) = f(k, x^k, s^k)$. Similarly, the concatenation of the encoder, the channel and the decoder, maps s^k and w^k to $u(k)$ as a possibly time-dependent function $u(k) = \psi(k, w^k, s^k)$. Under these assumptions, the following fundamental result was shown in [19] (Lemma 5.1):

$$I(r^k, p; u^k) \geq I(r^k; u^k) + I(p; e^k). \quad (8)$$

By further assuming in [19] that the decoder D in Figure 1b is deterministic, the following Markov chain naturally holds,

$$(p, r^k) \longleftrightarrow y^k \longleftrightarrow u^k, \quad (9)$$

leading directly to:

$$I(r^k, p; y^k) \geq I(r^k; u^k) + I(p; e^k), \quad (10)$$

which is found in the proof of [19] (Corollary 5.3). The deterministic nature of the decoder D played a crucial role in the proof of this result, since otherwise the Markov chain (9) does not hold, in general, due to the feedback from u to y .

Notice that both (8) and (10) provide lower bounds to mutual information as the sum of two mutual information terms, each of them relating a signal *external* to the loop (such as p, r^k) to a signal *internal* to the loop (such as u^k or y^k). Instead, the inequality:

$$I(x^k \rightarrow y^k) \geq I(r^k; y^k), \quad (11)$$

which holds for the system in Figure 1a and appears in [13] (Theorem 3) (and rediscovered later in [17] (Lemma 4.8.1)), involves the directed information between two internal signals and the mutual information between the second of these and an external sequence.

Remark 2. By using (22), $I(p^k, q^k, r^k; y^k) = I(r^k; y^k) + I(p^k, q^k; y^k | r^k)$. Then, applying Theorem 1, we recover (11), whenever $s \perp\!\!\!\perp (q, r, p)$. Thus, [13,17] (Theorem 3) and (Lemma 4.8.1)) can be obtained as a corollary of Theorem 1.

A related bound, similar to (10) but involving information rates and with the leftmost mutual information replaced by the directed information from x^k to y^k (which are two signals internal to the loop), has been obtained in [18] (Lemma 4.1) for the networked control system of Figure 1b:

$$\bar{I}(x \rightarrow y) \geq \bar{I}(r; u) + \lim_{k \rightarrow \infty} \frac{I(p; e^k)}{k}, \quad (12)$$

with $\bar{I}(x \rightarrow y) \triangleq \lim_{k \rightarrow \infty} \frac{1}{k} I(x^k \rightarrow y^k)$ and $\bar{I}(r; u) \triangleq \lim_{k \rightarrow \infty} \frac{1}{k} I(r^k; u^k)$, provided $\sup_{i \geq 0} E[x(i)^T x(i)] < \infty$. This result relies on three assumptions: (a) that the channel f is memory-less and satisfies a “conditional invertibility” property, (b) a finite-memory condition, and (c) a fading-memory condition, the latter two related to the decoder D (see Figure 1).

It is worth noting that, as defined in [18], these assumptions upon D exclude the use of side information by the decoder and/or the possibility of D being affected by random noise or having a random internal state that is non-observable (please see [18] for a detailed description of these assumptions).

Remark 3. In Section 4 we present Theorem 4, which yields (12) as a special case, but for the general system of Figure 1a and with no other assumption than mutual independence between r, p, q, s . Moreover, since with this independence condition Theorem 1 yields $I(x^k \rightarrow u^k) = I(r^k, p^k; u^k)$, the same happens with (8).

The inequality (11) has been extended in [16] (Theorem 1), for the case of discrete-valued random variables and assuming $s \perp (r, p, q)$, as the following identity (written in terms of the signals and setup shown in Figure 1a):

$$I(x^k \rightarrow y^k) = I(p^k, y^k) + I(x^k \rightarrow y^k | p^k). \quad (13)$$

Letting $q = s$ in Figure 1a and with the additional assumption that $(p, s) \perp q$, it was also shown in [16] (Theorem 1) that:

$$I(x^k \rightarrow y^k) = I(p^k; y^k) + I(q^{k-1}; y^k) + I(p^k; q^{k-1} | y^k), \quad (14)$$

for the cases in which $u(i) = y(i) + q(i)$ (i.e., when the concatenation of S_4 and S_1 corresponds to a summing node). In [16], (13) and (14) play important roles in characterizing the capacity of channels with noisy feedback.

To the best of our knowledge, (8), (10), (11)–(14) are the only results available in the literature that lower bound the difference between internal-to-internal directed information and external-to-internal mutual information. There exist even fewer published results in relation to inequalities between two directed information terms involving only signals internal to the loop. To the best of our knowledge, the only inequality of this type in the literature is the one found in the proof of Theorem 4.1 of [20]. The latter takes the form of a (conditional) data-processing inequality for directed information in closed-loop systems, and states that:

$$I(x^k \rightarrow y^k \parallel q^k) \geq I(x^k \rightarrow u^k), \quad (15)$$

provided: $q \perp (r, p)$ and if S_4 is such that y^i is a function of (u^i, q^i) (i.e., if S_4 is conditionally invertible) $\forall i$.

Inequality (15) plays a crucial role in [20], since it allows [20] (Thm. 4.1) to lower bound the average data rate across a digital error-free channel by a directed information. The setup considered in that theorem is shown in Figure 3, where \mathcal{F} is a plant, and \mathcal{E}, \mathcal{D} are the (source) encoder and decoder, respectively. In this figure, the variables that have been adapted to match those in Figure 4a (r, p, x correspond to disturbance, initial state, and plant output, respectively). Assuming $(r, p) \perp (q, s)$ and a conditionally invertible decoder, and letting $R(i)$ be the expected length (in bits) necessary for a binary representation of $y(i)$ given q^i , it states that $\frac{1}{k} \sum_{i=1}^k R(i) \geq \frac{1}{k} I(x^k \rightarrow u^k)$, $k = 1, 2, \dots$. This is a key result, because, combined with [20] (Equation (9)), it yields:

$$\frac{1}{k} I(x^k \rightarrow u^k) \leq \frac{1}{k} \sum_{i=1}^k R(i) \leq \frac{1}{k} I(x^k \rightarrow u^k) + 1 \quad [\text{bits/sample}], \quad k = 1, 2, \dots \quad (16)$$

This result highlights the operational meaning of the directed information as a lower bound (tight to within one bit) to the data rate of any given source code in a closed-loop system. This fact has been a crucial ingredient in characterizing the best rate performance achievable in Gaussian linear quadratic networked control [23,31], demonstrating the relevance of directed data-processing inequalities.

Unfortunately, as we will reveal in [26], the proof of [20] (Theorem 4.1) turns out to be invalid, since it relies upon [20] (Lemma 4.2), whose first claim does not hold. In [26] we use Theorem 2 to prove Theorem 4.1 of [20] without requiring a conditionally invertible decoder. This further illustrates the applicability of our results.

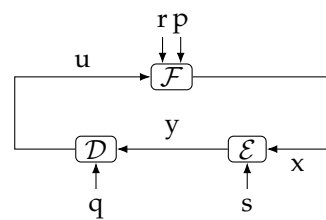


Figure 3. The networked control system considered in [20] (Figure 2), slightly simplified. The variables r, p, x, y correspond to d, x_o, y, s in [20], respectively.

In [23] (Lemma 1) another data-processing inequality is stated, which for the system in Figure 1a is equivalent to:

$$I(x^k \rightarrow y^k \| u_+^{k-1}) \geq I(x^k \rightarrow u^k), \quad k = 1, 2, \dots \quad (17)$$

where: $I(x^k \rightarrow y^k \| u_+^{k-1}) \triangleq \sum_{i=1}^k I(x^i; y(i) | y^{i-1}, u^{i-1})$. However, in [23] the blocks $\mathcal{S}_3, \mathcal{S}_4$ are defined implicitly, writing instead their input–output relation as collections of stochastic kernels $\mathbb{P}(y(i) | x^i, y^{i-1}), \mathbb{P}(u(i) | y^i, u^{i-1}), i = 1, 2, \dots$. The notation $\mathbb{P}(t|v, w)$ is to be understood as the conditional distribution of t given (v, w) . Crucially, this entails the implicit assumption that given x^i and y^{i-1} , $y(i)$ is independent of every other signal in the system (and likewise for $\mathbb{P}(u(i) | y^i, u^{i-1})$). In the representation of Figure 1a, this corresponds to assuming $q \perp\!\!\!\perp s$ and $(q, s) \perp\!\!\!\perp (r, p)$.

Remark 4. The conditioning on the side information q in both Theorem 2 and [20] (Theorem 4.1) is motivated by the use of entropy coded subtractively dithered quantization (ECSDQ) in obtaining the upper bound in (16). For such a scenario, the sequences q and s are identical and correspond to the dither signal, which is independent of r, p . This satisfies the requirements of (6) in Theorem 2 and of [20] (Theorem 4.1), but not the assumption that $q \perp\!\!\!\perp s$ and $(q, s) \perp\!\!\!\perp (r, p)$ implicit in [23] (Lemma 1), which yields (17). In spite of this, Lemma 1 of [23] is used in that paper to prove the lower bound in [23] (Equation (8)), an analogue of (16) which also considers the use of ECSDQ for the rate term and its upper bound.

1.3. Outline of the Paper

The remainder of the paper continues with some preliminary definitions and results in Section 2, the last of which is Lemma 2, a crucial tool for proving most of our theorems. Then follows the proof of Theorem 1 in Section 3. Section 4 presents additional inequalities relating mutual information between external–internal signal pairs, and directed information from one internal signal to another internal signal. These results can be seen as extensions or consequences of Theorem 1. Then we develop in Section 5 inequalities between two nested directed information expressions. Such results are the precursors of Theorem 2 and, as such, play a key role in its proof (which opens Appendix A). The notions and results associated with in-the-loop channel coding are developed in Section 6. The main conclusions of this work are presented in Section 7. Appendix A provides the proofs that are not written right after their corresponding theorems.

An earlier version of this work was made publicly available on arxiv.org [32] and, as such, it was cited in [23,31,33–35].

2. Preliminaries

2.1. Notation

The set of natural numbers is denoted \mathbb{N} . Random variables are denoted using non-italic characters, such as x . We write x^i to represent the sequence $x(1), x(2), \dots, x(i)$. We write $x \perp\!\!\!\perp y$ to express that x and y are independent. We use $\Pr\{\text{“outcome”}\}$ to denote the probability of a specific outcome of one or more random variables. For example, $\Pr\{x = 1, y \in \mathcal{Y}\}$ is the probability that $x = 1$ and y is in a given set \mathcal{Y} . Likewise, $\Pr\{\text{“outcome 1”} | \text{“outcome 2”}\}$ is the conditional probability of “outcome 1” given “out-

come 2". The Markov-chain notation $x \leftrightarrow y \leftrightarrow z$ means $\Pr\{x \in \mathcal{X}, z \in \mathcal{Z} | y \in \mathcal{Y}\} = \Pr\{x \in \mathcal{X} | y \in \mathcal{Y}\} \Pr\{z \in \mathcal{Z} | y \in \mathcal{Y}\}$, for every choice of the sets \mathcal{X} , \mathcal{Y} , \mathcal{Z} in the event spaces of x , y , and z , respectively. For two probability measures μ , ν on a common event space \mathcal{U} the notation $\mu \ll \nu$ means that μ is absolutely continuous with respect to ν , i.e., that $\forall \mathcal{U} \in \mathcal{U} : \nu(\mathcal{U}) = 0 \Rightarrow \mu(\mathcal{U}) = 0$.

2.2. Mutual Information

Let (Ω, \mathcal{F}, P) be a probability space, and $(\mathcal{X}, \mathcal{F}_{\mathcal{X}})$ and $(\mathcal{Y}, \mathcal{F}_{\mathcal{Y}})$ be measurable spaces, and consider the random variables $x : \Omega \rightarrow \mathcal{X}$, $y : \Omega \rightarrow \mathcal{Y}$. Define $\mathcal{M} \triangleq \mathcal{F}_{\mathcal{X}} \otimes \mathcal{F}_{\mathcal{Y}}$, i.e., the σ -algebra generated by the rectangles $\{A \times B : A \in \mathcal{F}_{\mathcal{X}}, B \in \mathcal{F}_{\mathcal{Y}}\}$. Consider a probability space $(\mathcal{X} \times \mathcal{Y}, \mathcal{M}, m)$ where m is the (joint) distribution of (x, y) , i.e., $m = P \circ (x, y)^{-1}$.

Denote the marginal probability distributions of x and y by μ and ν , respectively, where:

$$\mu(A) = m(A \times \mathcal{Y}), \quad A \in \mathcal{F}_{\mathcal{X}} \quad (18)$$

$$\nu(B) = m(\mathcal{X} \times B), \quad B \in \mathcal{F}_{\mathcal{Y}} \quad (19)$$

Define the product measure $\pi \triangleq \mu \times \nu$ on $(\mathcal{X} \times \mathcal{Y}, \mathcal{M})$.

Definition 1. With the above definitions, the mutual information between x and y is defined as:

$$I(x; y) \triangleq \int \log \left(\frac{dm}{d\pi} \right) d\pi, \quad (20)$$

where $\frac{dm}{d\pi}$ is the Radon–Nikodym derivative of m with respect to π [36].

An ensemble of random variables has a standard alphabet if it takes values from a set \mathcal{A} in a standard measurable space $(\mathcal{A}, \mathcal{F}_{\mathcal{A}})$ [37] [Section 1.4] and its probability measure is defined on $\mathcal{F}_{\mathcal{A}}$. For our purposes, it suffices to say that standard alphabets include discrete spaces, the real line, Euclidean vector spaces, and Polish spaces (i.e., complete separable metric spaces) [38].

Lemma 1 (Chain Rule of Mutual Information from [37] (Corollary 7.14)). Suppose x, y, z are random variables with standard alphabets and with joint distribution P_{xyz} . Suppose also that there exists a product distribution $M_{xyz} = M_x \times M_{yz}$ such that $M_{xyz} \gg P_{xyz}$. (This is true, for example, if $P_x \times P_{yz} \gg P_{xyz}$.) Then:

$$I(x; y, z) = I(x; y) + I(x; z | y). \quad (21)$$

From [37] (Lemma 7.4 and Equation 7.28), we have that $I(x; y, z) < \infty \Rightarrow P_x \times P_{yz} \gg P_{xyz}$, and thus (21) also holds if $I(x; y, z)$ is finite.

The conditional version of the chain rule of mutual information [39] (see also [37] (Corollary 2.5.1)) will be extensively utilized in the proofs of our results:

$$I(t, v; w | z) = I(v; w | z) + I(t; w | v, z). \quad (22)$$

For discrete random variables x, y , taking values from the sets \mathcal{X}, \mathcal{Y} , respectively, the entropy of x is defined as:

$$H(x) \triangleq - \sum_{x \in \mathcal{X}} \Pr\{x = x\} \log(\Pr\{x = x\}) \quad (23)$$

and the conditional entropy of x given y is defined as:

$$H(x | y) \triangleq - \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} \Pr\{x = x, y = y\} \log(\Pr\{x = x | y = y\}) \quad (24)$$

The entropy satisfies the chain rule:

$$H(x, y) = H(x) + H(y | x) = H(y) + H(x | y) \quad (25)$$

and is related to the mutual information as:

$$I(x; y) = H(x) + H(y) - H(x, y) = H(x) - H(x | y) = H(y) - H(y | x). \quad (26)$$

2.3. System Description

We begin by providing a formal description of the systems labeled $\mathcal{S}_1 \dots \mathcal{S}_4$ in Figure 1a. Their input–output relationships are given by the possibly-varying deterministic mappings (For notational simplicity, we omit writing their time dependency explicitly):

$$e(i) = \mathcal{S}_1(u^{i-d_1(i)}, r^i), \quad (27a)$$

$$x(i) = \mathcal{S}_2(e^{i-d_2(i)}, p^i), \quad (27b)$$

$$y(i) = \mathcal{S}_3(x^{i-d_3(i)}, s^i), \quad (27c)$$

$$u(i) = \mathcal{S}_4(y^{i-d_4(i)}, q^i), \quad (27d)$$

where r, p, s, q are exogenous random signals and the (possibly time-varying) delays $d_1, d_2, d_3, d_4 \in \{0, 1, \dots\}$ are such that:

$$d_1(k) + d_2(k) + d_3(k) + d_4(k) \geq 1, \quad \forall k \in \mathbb{N}.$$

That is, the concatenation of $\mathcal{S}_1, \dots, \mathcal{S}_4$ has a delay of at least one sample. For every $i \in \{1, \dots, k\}$, $r(i) \in \mathbb{R}^{n_r(i)}$, i.e., $r(i)$ is a real random vector whose dimension is given by some function $n_r : \{1, \dots, k\} \rightarrow \mathbb{N}$. The other sequences (q, p, s, x, y, u) are defined likewise.

2.4. A Necessary Modification of the Definition of Directed Information

As stated in [13], the directed information (as defined in (3)) is a more meaningful measure of the flow of information between x^k and y^k than the conventional mutual information $I(x^k; y^k) = \sum_{i=1}^k I(y(i); x^i | y^{i-1})$ when there exists causal feedback from y to x . In particular, if x^k and y^k are discrete-valued sequences, the input and output, respectively, of a forward channel, and if there exists *strictly causal* perfect feedback, so that $x(i) = y(i-1)$ (a scenario utilized in [13] as part of an argument in favor of the directed information), then the mutual information becomes:

$$\begin{aligned} I(x^k; y^k) &= H(y^k) - H(y^k | x^k) = H(y^k) - H(y^k | y^{k-1}) = H(y^k) - H(y(k) | y^{k-1}) \\ &= H(y^{k-1}). \end{aligned}$$

Thus, when strictly causal feedback is present, $I(x^k; y^k)$ fails to account for how much information about x^k has been conveyed to y^k through the forward channel that lies between them.

It is important to note that in [13] (as well as in many works concerned with communications), the forward channel is instantaneous, i.e., it has no delay. Therefore, if a feedback channel is utilized, then this feedback channel must have a delay of at least one sample, as in the example above. However, when studying the system in Figure 1a, we may need to evaluate the directed information between signals x^k and y^k which are, respectively, the input and output of a *strictly casual* forward channel (i.e., with a delay of at least one sample), whose output is instantaneously fed back to its input. In such a case, if one further assumes perfect feedback and sets $x(i) = y(i)$, then, in the same spirit as before,

$$I(x^k \rightarrow y^k) = \sum_{i=1}^k I(y(i); x^i | y^{i-1}) = \sum_{i=1}^k [H(y(i) | y^{i-1}) - H(y(i) | x^i, y^{i-1})] = H(y^k).$$

As one can see, Massey's definition of directed information ceases to be meaningful if instantaneous feedback is utilized.

It is natural to solve this problem by recalling that, in the latter example, the forward channel had a delay, say d , greater than one sample. Therefore, if we are interested in measuring how much of the information in $y(i)$, not present in y^{i-1} , was conveyed from x^i through the forward channel, we should look at the mutual information $I(y(i); x^{i-d} | y^{i-1})$, because only the input samples x^{i-d} can have an influence on $y(i)$. For this reason, we introduce the following, modified notion of directed information.

Definition 2 (Directed Information with Forward Delay). *In this paper, the directed information from x^k to y^k through a forward channel with a non-negative time varying delay of $d_{xy}(i)$ samples is defined as:*

$$I(x^k \rightarrow y^k) \triangleq \sum_{i=1}^k I(y(i); x^{i-d_{xy}(i)} | y^{i-1}). \quad (28)$$

For a zero-delay forward channel, the latter definition coincides with Massey's [13].

Likewise, we adapt the definition of causally-conditioned directed information to the definition:

$$I(x^k \rightarrow y^k \| e^k) \triangleq \sum_{i=1}^k I(y(i); x^{i-d_{xy}(i)} | y^{i-1}, e^i).$$

where, as before, $d_{xy}(i)$ is the delay from x to $y(i)$.

2.5. A Fundamental Lemma

The following result is an essential ingredient in the proof of most of our theorems:

Lemma 2. *In the system shown in Figure 4, the exogenous signals r, q are mutually independent and S_1, S_2 are deterministic (possibly time-varying) causal measurable functions characterized by $y^i = S_1(r^i, u^i)$, $u^i = S_2(q^i, y^{i-1})$, $\forall i \in \{1, \dots\}$, with $y_0 = y_0$ (deterministic). For this system, and for every $0 \leq j \leq i \leq k$ such that $i - j \leq 1$ and $i \geq 1$, the following Markov chain holds:*

$$r^k \longleftrightarrow (u^i, y^j) \longleftrightarrow q^k, \quad \forall k \in \mathbb{N}. \quad (29)$$

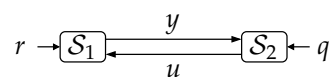


Figure 4. Two arbitrary causal systems S_1, S_2 interconnected in a feedback loop. The exogenous signals r, q are mutually independent.

Proof. Let $\mathcal{R}, \mathcal{Q}, \mathcal{U}, \mathcal{Y}$ be the event spaces of r^k, q^k, u^i, y^j , respectively. Since $y^j = S_1(r^j, u^j)$ and $u^i = S_2(q^i, y^{i-1})$ are deterministic measurable functions, it follows that for every possible pair of events $U \in \mathcal{U}, Y \in \mathcal{Y}$, the preimage sets $\mathcal{R}_{U,Y} \triangleq \{r^k : S_1(r^j, u^j) \in Y, u^i \in U\}$ and $\mathcal{Q}_{U,Y} \triangleq \{q^k : S_2(q^i, y^{i-1}) \in U, y^j \in Y\}$ are also deterministic and belong to \mathcal{R} and \mathcal{Q} , respectively. Thus, $(u^i, y^j) \in U \times Y \iff (r^k \in \mathcal{R}_{U,Y}, q^k \in \mathcal{Q}_{U,Y})$. This means that for every pair of events $R \in \mathcal{R}, Q \in \mathcal{Q}$,

$$\begin{aligned}
& \Pr\{r^k \in R, q^k \in Q | y^j \in \mathcal{Y}, u^i \in \mathcal{U}\} \\
& \stackrel{(a)}{=} \Pr\{r^k \in R, q^k \in Q | r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}, q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\} \\
& \stackrel{(b)}{=} \frac{\Pr\{r^k \in R \cap \mathcal{R}_{\mathcal{U}, \mathcal{Y}}, q^k \in Q \cap \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\}}{\Pr\{r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}, q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\}} \\
& \stackrel{(c)}{=} \frac{\Pr\{r^k \in R \cap \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\}}{\Pr\{r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\}} \cdot \frac{\Pr\{q^k \in Q \cap \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\}}{\Pr\{q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\}} \\
& = \frac{\Pr\{r^k \in R \cap \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\} \Pr\{q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\}}{\Pr\{r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\} \Pr\{q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\}} \cdot \frac{\Pr\{q^k \in Q \cap \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\} \Pr\{r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\}}{\Pr\{q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\} \Pr\{r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\}} \\
& \stackrel{(d)}{=} \frac{\Pr\{r^k \in R \cap \mathcal{R}_{\mathcal{U}, \mathcal{Y}}, q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\}}{\Pr\{r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}, q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\}} \cdot \frac{\Pr\{q^k \in Q \cap \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}, r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\}}{\Pr\{q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}, r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\}} \\
& \stackrel{(e)}{=} \Pr\{r^k \in R | r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}, q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}\} \cdot \Pr\{q^k \in Q | q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}}, r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}\} \\
& \stackrel{(f)}{=} \Pr\{r^k \in R | y^j \in \mathcal{Y}, u^i \in \mathcal{U}\} \cdot \Pr\{q^k \in Q | y^j \in \mathcal{Y}, u^i \in \mathcal{U}\}
\end{aligned}$$

where (a) and (f) follow because of the equivalence between the events $(y^j \in \mathcal{Y}, u^i \in \mathcal{U})$ and $(r^k \in \mathcal{R}_{\mathcal{U}, \mathcal{Y}}, q^k \in \mathcal{Q}_{\mathcal{U}, \mathcal{Y}})$, (b) and (e) follow from Bayes rule, and (c) and (d) are true because $r^k \perp q^k$. This completes the proof. \square

3. Proof of Theorem 1

It is clear from Figure 1a and from (27) that the relationship between r, p, q, s, x , and y can be represented by the diagram shown in Figure 5.

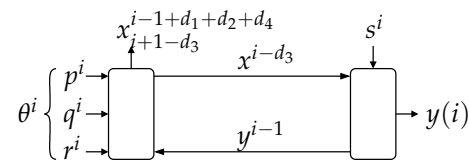


Figure 5. Representation of the system of Figure 1b highlighting the dependency between p, q, r, s, x , and y . The dependency on i of the delays $d_1(i), \dots, d_4(i)$ is omitted for clarity.

From this diagram and Lemma 2 it follows that if s is independent of (r, p, q) , then the following Markov chain holds:

$$y(i) \longleftrightarrow (x^{i-d_3(i)}, y^{i-1}) \longleftrightarrow (p^i, q^i, r^i). \quad (30)$$

Denoting the triad of exogenous signals p^k, q^k, r^k by:

$$\theta^k \triangleq (p^k, q^k, r^k), \quad (31)$$

we have the following:

$$\begin{aligned}
I(x^k \rightarrow y^i) &= \sum_{i=1}^k I(y(i); x^{i-d_3(i)} | y^{i-1}) \\
&\stackrel{(22)}{=} \sum_{i=1}^k \left[I(\theta^i, x^{i-d_3(i)}; y(i) | y^{i-1}) - I(\theta^i; y(i) | x^{i-d_3(i)}, y^{i-1}) \right] \\
&\stackrel{(a)}{=} \sum_{i=1}^k \left[I(\theta^i; y(i) | y^{i-1}) - I(\theta^i; y(i) | x^{i-d_3(i)}, y^{i-1}) \right] \tag{32a}
\end{aligned}$$

$$\stackrel{(b)}{\leq} \sum_{i=1}^k I(\theta^i; y(i) | y^{i-1}) \stackrel{(c)}{\leq} \sum_{i=1}^k I(\theta^k; y(i) | y^{i-1}) \tag{32b}$$

$$= I(\theta^k; y^k). \tag{32c}$$

In the above, (a) follows from the fact that, if y^{i-1} is known, and then $x^{i-d_3(i)}$ is a deterministic function of θ^i . The resulting sums on the right-hand side of (32a) correspond to $I(q^k, r^k, p^k \rightarrow y^k) - I(q^k, r^k, p^k \rightarrow y^k \parallel x^k)$, thereby proving the first part of the theorem, i.e., the equality in (5). In turn, (b) stems from the non-negativity of mutual information turn into equality if $s \perp\!\!\!\perp (r, p, q)$, as a direct consequence of the Markov chain in (30). Finally, equality holds in (c) if $s \perp\!\!\!\perp (q, r, p)$, since y depends causally upon θ . This shows that equality in (5) is achieved if $s \perp\!\!\!\perp (q, r, p)$, completing the proof.

4. Relationships between Mutual and Directed Information

The following result provides an inequality relating $I(x^k \rightarrow y^k)$ with the separate flows of information $I(r^k; y^k)$ and $I(p^k, q^k; y^k)$.

Theorem 3. For the system shown in Figure 1a, if $s \perp\!\!\!\perp (p, q, r)$ and $r^k \perp\!\!\!\perp (p^k, q^k)$, then:

$$I(x^k \rightarrow y^k) \geq I(r^k; y^k) + I(p^k, q^k; y^k). \tag{33}$$

with equality if and only if the Markov chain $(p^k, q^k) \leftrightarrow y^k \leftrightarrow r^k$ holds.

Theorem 3 shows that, provided $(p, q, r) \perp\!\!\!\perp s$, $I(x^k \rightarrow y^k)$ is lower bounded by the sum of the individual flows from all the subsets in any given partition of (p^k, q^k, r^k) , to y^k , provided these subsets are mutually independent. Indeed, both Theorems 1 and 3 can be generalized for any appropriate choice of external and internal signals. More precisely, let Θ be the set of all external signals in a feedback system. Let α and β be two internal signals in the loop. Define $\Theta_{\alpha, \beta} \subset \Theta$ as the set of exogenous signals that are introduced to the loop at every subsystem \mathcal{S}_i that lies in the path going from α to β . Thus, for any $\rho \in \Theta \setminus \Theta_{\alpha, \beta}$, if $\Theta_{\alpha, \beta} \perp\!\!\!\perp \Theta \setminus \Theta_{\alpha, \beta}$, we have that (5) and (33) become:

$$I(\alpha \rightarrow \beta) = I(\Theta \setminus \{\Theta_{\alpha, \beta}\}; \beta), \tag{34}$$

$$I(\alpha \rightarrow \beta) - I(\rho; \beta) \geq I(\Theta \setminus \{\rho \cup \Theta_{\alpha, \beta}\}; \beta), \tag{35}$$

respectively.

To finish this section, we present a stronger, non-asymptotic version of inequality (12):

Theorem 4. In the system shown in Figure 1a, if (r, p, q, s) are mutually independent, then:

$$I(x^k \rightarrow y^k) = I(r^k; u^k) + I(p^k; e^k) + I(q^k; y^k) + I(p^k; u^k | e^k) + I(r^k, p^k; y^k | u^k). \tag{36}$$

Remark 5. As anticipated, Theorem 4 can be seen as an extension of (12) to the more general setup shown in Figure 1a, where the assumptions made in [18] (Lemma 4.1) do not need to hold. In particular, letting the decoder D and p in Figure 1b correspond to S_4 and p^k in Figure 1a,

respectively, we see that inequality (12) holds even if the channel f has memory or D and E have independent initial states, or if the internal state of D is not observable [40].

Theorem 4 also admits an interpretation in terms of information flows. This can be appreciated in the diagram shown in Figure 6, which depicts the individual full-turn flows (around the entire feedback loop) stemming from q , r , and p . Theorem 4 states that the sum of these individual flows is a lower bound for the directed information from x to y , provided q, r, p, s are independent.

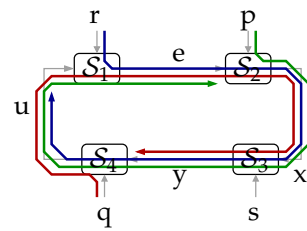


Figure 6. A representation of the three first information flows on the right-hand side of (36).

5. Relationships between Nested Directed Information

This section presents three closed-loop versions of the data-processing inequality relating two directed information terms, both between pairs of signals *internal* to the loop. As already mentioned in Section 1, to the best of our knowledge, the first inequality of this type to appear in the literature is the one in Theorem 4.1 in [20] (see (15)). Recall that the latter result stated that $I(x^k \rightarrow y^k \parallel q^k) \geq I(x^k \rightarrow u^k)$, requiring S_4 to be such that y^i is a deterministic function of (u^i, q^i) and that $q \perp (r, p)$. The following result presents another inequality that also relates two nested directed information terms, namely, $I(x^k \rightarrow y^k)$ and $I(e^k \rightarrow y^k)$, but requiring only that $s \perp (q, r, p)$.

Theorem 5. For the closed-loop system in Figure 1b, if $(q, r, p) \perp s$, then:

$$I(x^k \rightarrow y^k) \geq I(e^k \rightarrow y^k). \quad (37)$$

Notice that Theorem 5 does not require p to be independent of r or q . This may seem counterintuitive upon noting that p enters the loop between the link from e to x .

The following theorem is an identity between two directed information terms involving only internal signals. It can also be seen as a complement to Theorem 5, since it can be directly applied to establish the relationship between $I(e^k \rightarrow y^k)$ and $I(e^k \rightarrow u^k)$.

Theorem 6. For the system shown in Figure 1a, if $s \perp (q, r, p)$, then:

$$I(x^k \rightarrow y^k) \geq I(x^k \rightarrow u^k) + I(q^k; y^k) + I(r^k, p^k; y^k | u^k) + I(q^k; r^k | u^k, y^k), \quad (38)$$

with equality if, in addition, $q \perp (r, p)$. In the latter case, it holds that:

$$I(x^k \rightarrow y^k) = I(x^k \rightarrow u^k) + I(q^k; y^k) + I(r^k, p^k; y^k | u^k). \quad (39)$$

Notice that by requiring additional independence conditions upon the exogenous signals (specifically, $q \perp s$), Theorem 6 (and, in particular, (39)) yields:

$$I(x^k \rightarrow y^k) \geq I(x^k \rightarrow u^k), \quad (40)$$

which strengthens the inequality in [20] (Theorem 4.1) (stated above in (15)). More precisely, (40) does not require conditioning one of the directed information terms and holds irrespective of the invertibility of the mappings in the loop.

6. Giving Operational Meaning to the Directed Information: In-the-Loop Channel Coding

In this section we introduce the notions of in-the-loop transmission rate and capacity and show that they are related by the directed information rate across the channel in the same feedback loop. This provides another example to illustrate the applicability of Theorems 1 and 2 and also provides further operational meaning to the directed information rate.

Consider the scheme shown in Figure 7, and suppose \mathcal{C} is a noisy communication channel. Let \mathcal{E} and \mathcal{D} be the channel encoder and decoder, respectively, with r and p being side information sequences causally and independently available to each of them such that $(r, p) \perp (s, q)$. This means that, for $k = 1, 2, \dots, n$,

$$(p_1^n, r_{k+1}^n) \leftrightarrow r_1^k \leftrightarrow (w_1^{k+1}, x_1^k, y_0^k). \quad (41)$$

$$p_{k+1}^n \leftrightarrow p_1^k \leftrightarrow r_1^k \quad (42)$$

$$p_{k+1}^n \leftrightarrow p_1^k \leftrightarrow (w_1^{k+1}, x_1^k, y_0^k). \quad (43)$$

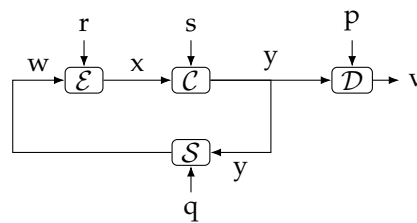


Figure 7. A communication feedback system in which the messages and the channel output are within the loop.

A crucial aspect of this scenario is the fact that the messages $w_1^n, w_{n+1}^{2n}, \dots$ to be encoded are contained in the sequence w , a signal internal to the loop; they can be regarded as a corrupted version of the decoded messages, which comprise the sequence v . This is a key difference with respect to the available literature on feedback capacity, where, to the best of the authors' knowledge, the messages are exogenous and the feedback signal only helps in the encoding task (exceptions can be found in some papers on networked control which consider in-the-loop channel coding, such as, e.g., [41,42]). In Figure 7, the latter standard scenario corresponds to encoding the sequence r .

The fact that the messages to be encoded bear information from the decoded message symbols requires one to redefine the notion of information transmission rate commonly used in the standard scenario. To see this, let $w(k) \in \mathcal{W}$, $k = 1, 2, \dots$, for some finite alphabet \mathcal{W} of cardinality $|\mathcal{W}|$, and notice that the transmission rate definitions $\log(|\mathcal{W}|)$ and $H(w_1^n)/n$ are unsatisfactory if $w(k) = y(k-1)$, $k = 1, 2, \dots$, i.e., if the messages to be transmitted are already available at the decoder (more generally, if there is no randomness in the feedback path). This suggests that a suitable notion of transmission rate for this scenario should exclude information that is already known by the receiver.

In view of the above, we propose the following notion of transmission rate for the case in which the messages to be transmitted are in the loop:

Definition 3. For the system described in Figure 1, the in-the-loop (ITL) transmission rate is defined as:

$$R_{ITL}^n \triangleq \frac{1}{n} \sum_{k=1}^n H(w(k) | w_1^{k-1}, y_0^{k-1}, p_1^k). \quad (44)$$

The meaning of the ITL transmission rate is further elucidated by considering the following scenarios:

1. If the feedback channel is deterministic, then $w(k)$ is a deterministic function of y_0^{k-1} and thus $R_{ITL}^n = 0$, as desired.
2. If the (forward) communication channel is noiseless, then at each time $k - 1$, we have $y_1^{k-1} = w_1^{k-1}$. Therefore $R_{ITL}^n = H(w_1^n | y_0^n, p_1^n) / n$. Again, if the feedback channel is deterministic, the ITL transmission rate is zero.
3. In the absence of feedback, $R_{ITL}^n = \frac{1}{n} H(w_1^n)$, recovering the notion of transmission rate of the case in which the messages are exogenous to the loop.

Thus, R_{ITL}^n can be interpreted as the sum of the information the encoder attempts to transmit at each sample time (expressed by the conditional entropy $H(w(k) | w_1^{k-1}, y_0^{k-1}, p_1^k)$) that is novel for the transmitter (because of the conditioning on w_1^{k-1}) and novel for the receiver (due to the conditioning on y_0^{k-1}, p_1^k).

Theorem 7. Consider the setup depicted in Figure 7, where \mathcal{E} and \mathcal{D} are the channel encoder and decoder, respectively, and \mathcal{C} is the communication channel. Suppose the message and side-information samples $w(k) \in \mathcal{W}$, $r(k) \in \mathcal{R}$, $k = 1, 2, \dots$, respectively, where \mathcal{W} and \mathcal{R} are finite alphabets. Define the binary random variable e_n to equal 1 if $v_1^n \neq w_1^n$ and 0 otherwise. Then, for every $n \in \mathbb{N}$,

$$R_{ITL}^n \geq I(w_1^n \rightarrow y_0^n | p_1^n), \quad (45)$$

with equality if and only if $H(w_1^n | y_0^n, p_1^n) = 0$. Moreover,

$$\Pr\{e_n = 1\} = \frac{R_{ITL}^n - \frac{1}{n} I(w_1^n \rightarrow y_0^n | p_1^n) - \frac{1}{n} H(e_n | y_0^n, p_1^n)}{\frac{1}{n} H(w_1^n | y_0^n, p_1^n, e_n = 1)} \quad (46)$$

$$\geq \frac{R_{ITL}^n - \frac{1}{n} I(w_1^n \rightarrow y_0^n | p_1^n) - 1/n}{\log_2(|\mathcal{W}|)} \quad (47)$$

Proof. Recall that:

$$I(w_1^n \rightarrow y_0^n | p_1^n) = \sum_{k=1}^n I(w_1^k, y(k) | y_0^{k-1}, p_1^k) = \sum_{k=1}^n H(w_1^k | y_0^{k-1}, p_1^k) - \sum_{k=1}^n H(w_1^k | y_0^k, p_1^k) \quad (48)$$

On the other hand,

$$nR_{ITL}^n = \sum_{k=1}^n H(w(k) | w_1^{k-1}, y_0^{k-1}, p_1^k) \stackrel{(cr)}{=} \sum_{k=1}^n H(w_1^k | y_0^{k-1}, p_1^k) - \sum_{k=2}^n H(w_1^{k-1} | y_0^{k-1}, p_1^k) \quad (49)$$

$$\stackrel{(48)}{=} \sum_{k=1}^n H(w_1^k | y_0^k, p_1^k) - \sum_{k=2}^n H(w_1^{k-1} | y_0^{k-1}, p_1^k) + I(w_1^n \rightarrow y_0^n | p_1^n) \quad (50)$$

$$\stackrel{(43)}{=} H(w_1^n | y_0^n, p_1^n) + I(w_1^n \rightarrow y_0^n | p_1^n), \quad (51)$$

where the equality (cr) follows from the chain rule of entropy. This proves the first part of the theorem.

Let us now re-derive the first steps leading to Fano's inequality, to include the side-information p_1^n and to verify that it is not affected by the fact that w and y are within the loop.

$$H(w_1^n | y_0^n, p_1^n) \stackrel{(cr)}{=} H(w_1^n, e_n | y_0^n, p_1^n) - H(e_n | y_0^n, p_1^n, w_1^n) \quad (52)$$

$$\stackrel{(a)}{=} H(e_n | y_0^n, p_1^n) + H(w_1^n | y_0^n, p_1^n, e_n) \quad (53)$$

$$\stackrel{(b)}{=} H(e_n | y_0^n, p_1^n) + H(w_1^n | y_0^n, p_1^n, e_n = 1) \Pr\{e_n = 1\}, \quad (54)$$

where the equality (cr) follows from the chain rule of entropy and (a) holds because $H(e_n | y_0^n, p_1^n, w_1^n) = 0$ and from the chain rule, while (b) is because $H(w_1^n | y_0^n, p_1^n, e_n = 0) = 0$.

Substituting this into (51),

$$nR_{ITL}^n = H(e_n | y_0^n, p_1^n) + H(w_1^n | y_0^n, p_1^n, e_n = 1) \Pr\{e_n = 1\} + I(w_1^n \rightarrow y_0^n \| p_1^n). \quad (55)$$

Noting that $H(e_n | y_0^n, p_1^n) \leq 1$ and $H(w_1^n | y_0^n, p_1^n, e_n = 1) \leq n \log(|\mathcal{W}|)$ leads directly to (46), comparing the proof. \square

Theorem 7 allows one to draw an additional interpretation of the ITL transmission rate. We extend first the identity of [27] to include causal conditioning by p_1^n :

$$I(w_1^n; y_0^n, p_1^n) = I(w_1^n; y_n | y_0^{n-1}, p_1^n) + I(w_n; y_0^{n-1} | w_1^{n-1}, p_1^n) + I(w_1^{n-1}; y_0^{n-1}, p_1^{n-1}) \quad (56)$$

$$= \sum_{k=1}^n I(w_1^k; y(k) | y_0^{k-1}, p_1^k) + \sum_{k=1}^n I(w(k); y_0^{k-1} | w_1^{k-1}, p_1^k) = I(w_1^n \rightarrow y_0^n \| p_1^n) + I(y_0^{n-1} \rightarrow w_1^n \| p_1^n), \quad (57)$$

where:

$$I(y_0^{n-1} \rightarrow w_1^n \| p_1^n) \triangleq \sum_{k=1}^n I(w(k); y_0^{k-1} | w_1^{k-1}, p_1^k). \quad (58)$$

It readily follows from (56) that:

$$H(w_1^n) - I(y_0^{n-1} \rightarrow w_1^n \| p_1^n) = I(w_1^n \rightarrow y_0^n \| p_1^n) + H(w_1^n | y_0^n, p_1^n) \stackrel{(51)}{=} nR_{ITL}^n. \quad (59)$$

Thus, the ITL transmission rate corresponds to the entropy rate of the messages having extracted from it the information flowing from the decoder input to the messages.

The main result of this section is the following theorem, which asserts that the supremum of achievable ITL transmission rates is upper bounded by the directed information across the communication channel.

Theorem 8. Consider the setup depicted in Figure 7, where \mathcal{E} and \mathcal{D} are the channel encoder and decoder, respectively, and \mathcal{C} is the communication channel. Then the supremum of achievable ITL transmission rates is upper bounded by the supremum of the directed information rate from x to y causally conditioned by p_1^n .

Proof. The result follows directly from Theorems 2 and 7. \square

Thus, the supremum of $\lim_{n \rightarrow \infty} I(x_1^n \rightarrow y_1^n \| p_1^n)$ is an outer bound to the capacity region of ITL transmission rates.

In the following example, this bound is reachable.

Example 1. Consider the case in which the forward channel \mathcal{C} in Figure 7 is transparent, i.e., $y(k) = x(k)$ for $k = 0, 1, \dots$, as shown in Figure 8. Let $y(k) \in \{0, 1, 2, 3\}$, $k = 0, 1, \dots$. Let $q(0) = 1$ (deterministically) and $q(1), q(2), \dots$ be binary and i.i.d. with $\Pr\{q(k) = 1\} = \alpha = 0.9$. The feedback channel \mathcal{S} is defined by the following recursion:

$$w(k) = \begin{cases} q(k) & , \text{ if } q(k-1) = (y(k-1) \bmod 2) \\ (y(k-1) \bmod 2) & , \text{ if } q(k-1) \neq (y(k-1) \bmod 2) \end{cases} , k = 1, 2, \dots \quad (60)$$

Thus, \mathcal{S} outputs a new sample of q iff the previous sample of q is matched by the previous sample $\bmod 2$ of y . Otherwise, it lets $y(k-1) \bmod 2$ pass through.

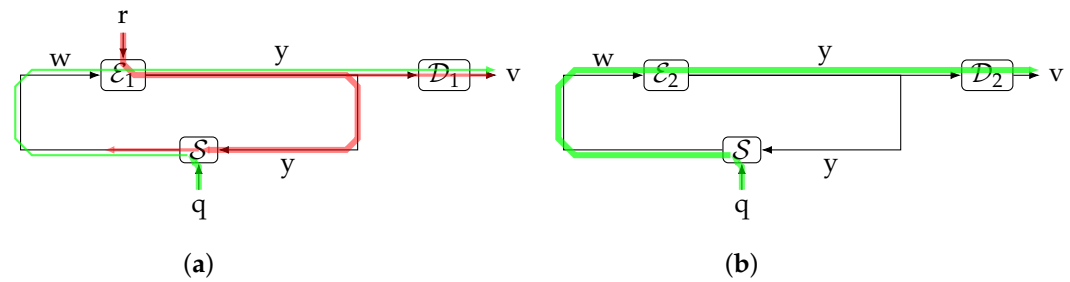


Figure 8. The feedback communication system considered in Example 1. In (a), encoder/decoder pair 1 yields a large mutual information between w and v by adding a backward information flow (in red) to the forward information flow (in green). The latter is only part of the entropy rate of q . In (b), encoder/decoder pair 2 yields a smaller mutual information between w and v , but it corresponds to the greatest possible forward information flow, which coincides with the entropy rate of q . Thus, it is capacity achieving with respect to the ITL transmission rate.

Consider first the following encoder–decoder pair, designed with the aim of achieving zero-error communication while maximizing $H(w_1^n)/n = I(w_1^n; v_1^n)$.

Encoder \mathcal{E}_1 : Let the side-information sequence r be binary i.i.d. and independent of q , with $\Pr\{r(k) = 1\} = \beta$, and:

$$y(0) = r(0) \quad (61)$$

$$y(k) = \begin{cases} r(k) & , \text{ if } w(k) = (y(k-1) \bmod 2) \\ r(k) + 2 & , \text{ if } w(k) \neq (y(k-1) \bmod 2) \end{cases} \quad , k = 1, 2, \dots \quad (62)$$

Decoder \mathcal{D}_1 :

$$v(k) = \begin{cases} y(k-1) \bmod 2 & , \text{ if } y(k) \leq 1 \\ (y(k-1) \bmod 2) \oplus 1 & , \text{ if } y(k) > 1 \end{cases} \quad (63)$$

where \oplus is the exclusive-or binary operator. With this choice, $v(k) = w(k)$ for $k = 1, 2, \dots$. In addition,

$$w(k) = \begin{cases} q(k) & , \text{ if } r(k-1) = q(k-1) \\ r(k-1) & , \text{ if } r(k-1) \neq q(k-1) \end{cases} \quad , k = 1, 2, \dots \quad (64)$$

Therefore,

$$\Pr\{w(1) = 1\} = \alpha\beta. \quad (65)$$

and, for $k \geq 2$,

$$\Pr\{w(k) = 1\} = (\alpha^2 + (1 - \alpha)^2)\beta + \alpha(1 - \alpha) \quad (66)$$

Thus, and since $\alpha = 0.9$, the entropy of each $w(k)$ is maximized by $\beta = 0.5055$. However, encoder \mathcal{E}_1 makes the samples of w interdependent, so finding the value of β that maximizes $H(w_1^n)/n$ (and thus $I(w_1^n; v_1^n)$ as well) is more involved, and that value does not need to be the same. We have found numerically that (for $n = 22$) the maximum of $H(w_1^n)/n = I(w_1^n; v_1^n)/n$ is (approximately) 0.9941 [bits/sample], attained with $\beta = 0.503$, very close to the β which maximizes $H(w_1^n)/n$.

For later comparison, we also calculate the value of R_{ITL}^n yielded by this choice of encoder:

$$R_{ITL}^n \stackrel{(a)}{=} I(w_1^n \rightarrow y_1^n) \stackrel{Thm. 1}{=} I(q_0^n; y_0^n) = \sum_{k=0}^n (H(q(k)|q_0^{k-1}) - H(q(k)|q_0^{k-1}, y_0^n)), \quad (67)$$

where (a) holds from Theorem 7 because $H(w_1^n | y_0^n) = 0$. Defining the binary random variables $t(k) \triangleq 1$ when $(y(k) \bmod 2) = q(k)$ and 0 otherwise, we obtain:

$$H(q(k) | q_0^{k-1}, y_0^n) = H(q(k) | q_0^{k-1}, y_0^n, t(k-1)) \quad (68)$$

$$\begin{aligned} &= H(q(k) | q_0^{k-1}, y_0^n, t(k-1) = 0) \Pr\{t(k-1) = 0\} + H(q(k) | q_0^{k-1}, y_0^n, t(k-1) = 1) \Pr\{t(k-1) = 1\} \\ &\stackrel{(64)}{=} H(q(k)) \Pr\{t(k-1) = 0\} + 0 \cdot \Pr\{t(k-1) = 1\} \end{aligned} \quad (69)$$

Thus,

$$R_{ITL}^n = H(q(k))(1 - \Pr\{t(k-1) = 0\}) = H(q(k))(\alpha(1 - \beta) + (1 - \alpha)\beta) \quad (70)$$

$$= 0.469 \times 0.4976 = 0.2334 \quad [\text{bits/sample}], \quad (71)$$

using $\beta = 0.503$.

The second encoder/decoder pair is set to maximize R_{ITL}^n , and is defined as follows:

Encoder \mathcal{E}_2 :

$$y(k) = \begin{cases} 1 & , \text{ if } k = 0 \\ w(k) & , \text{ if } k \geq 1 \end{cases} \quad (72)$$

Thus, zero-error communication is trivially attained with the simple decoding rule:

Decoder \mathcal{D}_2 :

$$v(k) = y(k), \quad k \geq 1. \quad (73)$$

In addition, encoder \mathcal{E}_2 yields $w(k) = q(k)$, for $k \geq 1$. Therefore,

$$\frac{1}{n} I(w_1^n \rightarrow y_0^n) \stackrel{\text{Thm. 7}}{=} R_{ITL}^n = \frac{1}{n} H(q_1^n) = 0.469 \quad [\text{bits/sample}] \quad (74)$$

As expected, encoder \mathcal{E}_2 yields a higher R_{ITL}^n than encoder \mathcal{E}_1 . More significant is the fact that encoder/decoder pair 2 achieves the in-the-loop capacity for this channel, since:

$$\frac{1}{n} I(w_1^n \rightarrow y_0^n) \stackrel{(a)}{=} I(q_1^n; y_0^n) \leq H(q_1^n) \quad (75)$$

The previous example illustrates an important fact that is closely related with the motivation behind the definition of R_{ITL}^n : maximizing the mutual information between the messages to be transmitted and the decoded messages (a leitmotif in traditional channel coding, wherein messages are generated outside the loop) is not suitable when messages are in the loop.

Indeed, (56) provides a mathematically precise meaning to the above observation. It reveals why maximizing $I(w_1^n; y_0^n, p_1^n)$ does not necessarily mean maximizing $I(w_1^n \rightarrow y_0^n \| p_1^n)$, since the former is the sum of backward and forward information flows (represented in green and red in Figure 8, respectively).

Finally, Theorems 7 and 8 imply that in the design of any encoder for in-the-loop messages, aiming to yield the joint probability distribution of channel input and output sequences that maximizes the directed information is of practical importance: it is necessary for achieving the highest “useful” transmission rate while minimizing the probability of error.

7. Concluding Remarks

The widely used data processing inequality does not hold for systems with feedback. In this work, we provided a very general *directed information* data processing inequality that is applicable to feedback systems. A key insight to be gained from this new inequality is that, for nested pairs of sequences, the further apart the signals in the feedback system are

from each other, the lower is the directed information between them (measuring distance from starting to finishing sequence and in the direction of cause and effect). Thus, post processing signals within a feedback loop cannot increase the information, which is similar to the open loop case. In order to obtain these results, we considered arbitrary causal systems that are interconnected in a feedback loop, with arbitrarily distributed signals. We were able to overcome the generally non-trivial dependencies between the signals in such a scenario by establishing a family of useful Markov chains that conditionally decouple the sequences in the system. These Markov chains are useful by themselves for studies involving interconnected systems. We further used the Markov chains to derive a number of fundamental information inequalities that are applicable to signals that are entirely within feedback loops or where some signals are inside and others outside the loop. With the use of these inequalities, we were able to show that the conventional notion of channel capacity is not adequate for *in-the-loop* communications. Instead, we provided the new notion of in-the-loop channel capacity, and described a special case where it was achievable. As an additional application of our results, we discussed how they allow one to generalize two known fundamental inequalities in networked control involving directed information. We are confident that our analysis provides useful insights to understand and think about information flows in single-loop feedback systems, and that our results will serve as a toolbox for research in, e.g., networked control systems or communications within a feedback loop.

There are several future research directions stemming from this work, from which we outline the following three:

1. Establishing whether (and under which conditions, if any) in the system of Figure 1, each of the following inequalities is true or false:

$$I(r^k; e^k) \geq I(r^k; x^k) \geq I(r^k; y^k) \geq I(r^k; u^k). \quad (76)$$

2. Extending Theorems 1 and 2 to scenarios with more than one feedback loop.
3. Exploring if tree codes [43] can be tailored to maximize the ITL data rate instead of the conventional data rate within a feedback loop. If such adaptation is possible, it would be interesting to assess how close to the ITL channel capacity such codes can perform.

Author Contributions: Both authors contributed to the derivation of the results presented in this paper and on its writing. All authors have read and agreed to the published version of the manuscript.

Funding: This research was partially funded by ANID-CONICYT Basal fund FB0008.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

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

Appendix A. Proofs

Proof of Theorem 2. If $(q, s) \perp\!\!\!\perp (r, p)$ and $q \perp\!\!\!\perp s$, then (6) follows by applying Theorem 5 and then Theorem 6. If $(p, s) \perp\!\!\!\perp (r, q)$ and $p \perp\!\!\!\perp s$, then one arrives at (6) by applying Theorem 6 followed by Theorem 5.

To prove the second part, notice that:

$$I(x^k \rightarrow y^k \parallel q^k) = I(x^k \rightarrow y^k | q^k) \quad (A1)$$

which follows since $x^{i-d_3(i)}, y^i$ are deterministic functions of (r^k, p^k, s^i, q^i) and $q_{i+1}^k \leftrightarrow q^i \leftrightarrow (r^k, p^k, s^i)$, a Markov chain that results from combining $q_{i+1}^k \leftrightarrow q^i \leftrightarrow s^i$ with $(q^k, s^k) \perp\!\!\!\perp (r^k, p^k)$.

On the other hand, the fact that $(r^k, p^k) \perp\!\!\!\perp (q^k, s^k)$ allows one to obtain from Theorem 1 that:

$$I(x^k \rightarrow y^k | q^k) = I(r^k, p^k; y^k | q^k). \quad (A2)$$

But:

$$\begin{aligned} I(r^k, p^k; y^k | q^k) &\stackrel{(22)}{=} I(r^k, p^k; u^k, y^k | q^k) - I(r^k, p^k; u^k | q^k, y^k) \\ &\stackrel{(a)}{=} I(r^k, p^k; u^k, y^k | q^k) \\ &\stackrel{(22)}{=} I(r^k, p^k; u^k, y^k, q^k) - I(r^k, p^k; q^k) \\ &\stackrel{(b)}{=} I(r^k, p^k; u^k, y^k, q^k) \\ &\stackrel{(22)}{=} I(r^k, p^k; u^k) + I(r^k, p^k; y^k, q^k | u^k) \end{aligned} \quad (A3)$$

where (a) is due to the fact that u^k is a deterministic function of q^k, y^k . Equality (b) holds if and only if $(r, p) \perp\!\!\!\perp q$. The fact that $(r^k, p^k) \perp\!\!\!\perp (q^k, s^k)$ allows one to obtain from Theorem 1 that $I(x^k \rightarrow u^k) = I(r^k, p^k; u^k)$. Substituting this into (A3) and then into (A2) and the latter into (A1), we obtain $I(x^k \rightarrow y^k | q^k) \geq I(x^k \rightarrow u^k)$, which combined with Theorem 5 yields (7). This completes the proof. \square

Proof of Theorem 3. Apply the chain-rule identity (22) to the *right-hand side* (RHS) of (5) to obtain:

$$I(\theta^k; y^k) = I(p^k, q^k, r^k; y^k) = I(p^k, q^k; y^k | r^k) + I(r^k; y^k). \quad (A4)$$

Now, applying (22) twice, one can express the term $I(p^k, q^k; y^k | r^k)$ as follows:

$$\begin{aligned} I(p^k, q^k; y^k | r^k) &= I(p^k, q^k; y^k, r^k) - I(p^k, q^k; r^k) = I(p^k, q^k; y^k, r^k) \\ &= I(p^k, q^k; y^k) + I(p^k, q^k; r^k | y^k), \end{aligned} \quad (A5)$$

where the second equality follows since $(p^k, q^k) \perp\!\!\!\perp r^k$. The result then follows directly by combining (A5) with (A4) and (5). \square

Proof of Theorem 4. Since $q \perp\!\!\!\perp (r, p, s)$,

$$I(x^k \rightarrow y^k) \stackrel{(a)}{=} I(x^k \rightarrow u^k) + I(q^k; y^k) + I(r^k, p^k; y^k | u^k) \quad (A6)$$

$$\stackrel{(b)}{=} I(r^k, p^k; u^k) + I(q^k; y^k) + I(r^k, p^k; y^k | u^k) \quad (A7)$$

$$\stackrel{(c)}{=} I(r^k; u^k) + I(p^k; u^k | r^k) + I(q^k; y^k) + I(r^k, p^k; y^k | u^k), \quad (A8)$$

where (a) is due to Theorem 6, (b) follows from Theorem 1 and the fact that $(s, q) \perp\!\!\!\perp (r, p)$, and (c) follows from the chain rule of mutual information. For the second term on the RHS of the last equation, we have:

$$I(p^k; u^k | r^k) \stackrel{(a)}{=} I(p^k; u^k | r^k) + I(p^k; r^k) = I(p^k; r^k, u^k) \quad (A9)$$

$$\stackrel{(b)}{=} I(p^k; r^k, u^k, e^k) - I(p^k; e^k | r^k, u^k) \quad (A10)$$

$$\stackrel{(c)}{=} I(p^k; r^k, u^k, e^k) \quad (A11)$$

$$\stackrel{(d)}{=} I(p^k; e^k) + I(p^k; r^k, u^k | e^k) \quad (A12)$$

$$\stackrel{(e)}{=} I(p^k; e^k) + I(p^k; u^k | e^k) + I(p^k; r^k | u^k, e^k) \quad (A13)$$

$$\stackrel{(f)}{=} I(p^k; e^k) + I(p^k; u^k | e^k), \quad (A14)$$

where (a) holds since $r \perp p$, (b), (d), and (e) stem from the chain rule of mutual information (22), and (c) is a consequence of the fact that $e^k = \mathcal{S}_1(u^{k-d_1(k)}, r^k)$. Finally, (f) is due to the Markov chain $r^k \leftrightarrow (u^k, e^k) \leftrightarrow p^k$, which holds because $r \perp (p, s, q)$ as a consequence of Lemma 2 in the Appendix (see also Figure 1a). Substitution of (A14) into (A8) yields (36), thereby completing the proof. \square

Proof of Theorem 5. Since $(p, q, r) \perp s$, we can apply (11) (where now (q, r) plays the role of r), and obtain

$$I(x^k \rightarrow y^k) \geq I(q^k, r^k; y^k). \quad (A15)$$

Now, we apply Theorem 1, which gives

$$I(q^k, r^k; y^k) \geq I(e^k \rightarrow y^k), \quad (A16)$$

completing the proof. \square

Proof of Theorem 6. We have that:

$$I(x^k \rightarrow y^k) \stackrel{(a)}{=} I(r^k, p^k, q^k; y^k) \stackrel{(22)}{=} I(q^k; y^k) + I(r^k, p^k; y^k | q^k) \quad (A17)$$

$$\stackrel{(A3)}{=} I(r^k, p^k; u^k) + I(r^k, p^k; y^k, q^k | u^k) \quad (A18)$$

$$\stackrel{(22)}{=} I(q^k; y^k) + I(r^k, p^k; u^k) + I(r^k, p^k; y^k | u^k) + I(r^k, p^k; q^k | u^k, y^k) \quad (A19)$$

$$\stackrel{(b)}{\geq} I(q^k; y^k) + I(x^k \rightarrow u^k) + I(r^k, p^k; y^k | u^k) + I(r^k, p^k; q^k | u^k, y^k) \quad (A20)$$

$$\stackrel{(c)}{\geq} I(q^k; y^k) + I(x^k \rightarrow u^k) + I(r^k, p^k; y^k | u^k),$$

where (a) follows from Theorem 1 and the assumption $(r, p, q) \perp s$, (b) is from Theorem 1, with equality iff $(q, s) \perp (r, p)$, and from Lemma 2 (in the Appendix), (c) turns into equality if $q \perp (r, p, s)$. This completes the proof. \square

References

1. Cover, T.M.; Thomas, J.A. *Elements of Information Theory*, 2nd ed.; Wiley-Interscience: Hoboken, NJ, USA, 2006.
2. Salek, S.; Cadamuro, D.; Kammerlander, P.; Wiesner, K. Quantum Rate-Distortion Coding of Relevant Information. *IEEE Trans. Inf. Theory* **2019**, *65*, 2603–2613. [\[CrossRef\]](#)
3. Lindenstrauss, E.; Tsukamoto, M. From Rate Distortion Theory to Metric Mean Dimension: Variational Principle. *IEEE Trans. Inf. Theory* **2018**, *64*, 3590–3609. [\[CrossRef\]](#)
4. Yang, Y.; Grover, P.; Kar, S. Rate Distortion for Lossy In-Network Linear Function Computation and Consensus: Distortion Accumulation and Sequential Reverse Water-Filling. *IEEE Trans. Inf. Theory* **2017**, *63*, 5179–5206. [\[CrossRef\]](#)
5. Derpich, M.S.; Østergaard, J. Improved upper bounds to the causal quadratic rate-distortion function for Gaussian stationary sources. *IEEE Trans. Inf. Theory* **2012**, *58*, 3131–3152. [\[CrossRef\]](#)

6. Ramakrishnan, N.; Iten, R.; Scholz, V.B.; Berta, M. Computing Quantum Channel Capacities. *IEEE Trans. Inf. Theory* **2021**, *67*, 946–960. [\[CrossRef\]](#)
7. Song, J.; Zhang, Q.; Kadhe, S.; Bakshi, M.; Jaggi, S. Stealthy Communication Over Adversarially Jammed Multipath Networks. *IEEE Trans. Inf. Theory* **2020**, *68*, 7473–7484. [\[CrossRef\]](#)
8. Makur, A. Coding Theorems for Noisy Permutation Channels. *IEEE Trans. Inf. Theory* **2020**, *66*, 6723–6748. [\[CrossRef\]](#)
9. Kostina, V.; Verdú, S. Lossy joint source-channel coding in the finite blocklength regime. *IEEE Trans. Inf. Theory* **2013**, *59*, 2545–2575. [\[CrossRef\]](#)
10. Huang, Y.; Narayanan, K.R. Joint Source-Channel Coding with Correlated Interference. *IEEE Trans. Commun.* **2012**, *60*, 1315–1327. [\[CrossRef\]](#)
11. Steinberg, Y.; Merhav, N. On hierarchical joint source-channel coding with degraded side information. *IEEE Trans. Inf. Theory* **2006**, *52*, 886–903. [\[CrossRef\]](#)
12. Hollands, S. Trace- and improved data processing inequalities for von Neumann algebras. *arXiv* **2021**, arXiv:2102.07479.
13. Massey, J.L. Causality, feedback and directed information. In Proceedings of the International Symposium on Information Theory and Its Applications, Honolulu, HI, USA, 27–30 November 1990; pp. 303–305.
14. Kramer, G. Directed Information for Channels with Feedback. Ph.D. Thesis, Swiss Federal Institute of Technology, Zürich, Switzerland, 1998.
15. Tatikonda, S.; Mitter, S. The Capacity of Channels With Feedback. *IEEE Trans. Inf. Theory* **2009**, *55*, 323–349. [\[CrossRef\]](#)
16. Li, C.; Elia, N. The Information Flow and Capacity of Channels with Noisy Feedback. *arXiv* **2011**, arXiv:1108.2815.
17. Tatikonda, S.C. Control Under Communication Constraints. Ph.D. Thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA, USA, 2000.
18. Martins, N.C.; Dahleh, M.A. Fundamental limitations of performance in the presence of finite capacity feedback. In Proceedings of the 2005 American Control Conference, Portland, OR, USA, 8–10 June 2005.
19. Martins, N.; Dahleh, M. Feedback control in the presence of noisy Channels: “Bode-like” fundamental limitations of performance. *IEEE Trans. Autom. Control* **2008**, *53*, 1604–1615. [\[CrossRef\]](#)
20. Silva, E.I.; Derpich, M.S.; Østergaard, J. A framework for control system design subject to average data-rate constraints. *IEEE Trans. Autom. Control* **2011**, *56*, 1886–1899. [\[CrossRef\]](#)
21. Silva, E.I.; Derpich, M.S.; Østergaard, J. On the Minimal Average Data-Rate That Guarantees a Given Closed Loop Performance Level. In Proceedings of the 2nd IFAC Workshop on Distributed Estimation and Control in Networked Systems, NECSYS, Annecy, France, 13–14 September 2010; pp. 67–72.
22. Silva, E.I.; Derpich, M.S.; Østergaard, J. An achievable data-rate region subject to a stationary performance constraint for LTI plants. *IEEE Trans. Autom. Control* **2011**, *56*, 1968–1973. [\[CrossRef\]](#)
23. Tanaka, T.; Esfahani, P.M.; Mitter, S.K. LQG Control With Minimum Directed Information: Semidefinite Programming Approach. *IEEE Trans. Autom. Control* **2018**, *63*, 37–52. [\[CrossRef\]](#)
24. Quinn, C.; Coleman, T.; Kiyavash, N.; Hatsopoulos, N. Estimating the directed information to infer causal relationships in ensemble neural spike train recordings. *J. Comput. Neurosci.* **2011**, *30*, 17–44. [\[CrossRef\]](#)
25. Permuter, H.H.; Kim, Y.H.; Weissman, T. Interpretations of directed information in portfolio theory, data Compression, and hypothesis testing. *IEEE Trans. Inf. Theory* **2011**, *57*, 3248–3259. [\[CrossRef\]](#)
26. Derpich, M.S.; Østergaard, J. Comments on “A Framework for Control System Design Subject to Average Data-Rate Constraints”. *arXiv* **2021**, arXiv:2103.12897.
27. Massey, J.; Massey, P. Conservation of mutual and directed information. In Proceedings of the Proceedings. International Symposium on Information Theory, 2005. ISIT 2005, Adelaide, SA, Australia, 4–9 September 2005; pp. 157–158. [\[CrossRef\]](#)
28. Kim, Y.H.; Kim, Y.H. A Coding Theorem for a Class of Stationary Channels With Feedback. *IEEE Trans. Inf. Theory* **2008**, *54*, 1488–1499. [\[CrossRef\]](#)
29. Zamir, R.; Kochman, Y.; Erez, U. Achieving the Gaussian rate-distortion function by prediction. *IEEE Trans. Inf. Theory* **2008**, *54*, 3354–3364. [\[CrossRef\]](#)
30. Zhang, H.; Sun, Y.X. Directed information and mutual information in linear feedback tracking systems. In Proceedings of the 6-th World Congress on Intelligent Control and Automation, Dalian, China, 21–23 June 2006; pp. 723–727.
31. Silva, E.I.; Derpich, M.S.; Østergaard, J.; Encina, M.A. A characterization of the minimal average data rate that guarantees a given closed-loop performance level. *IEEE Trans. Autom. Control* **2016**, *61*, 2171–2186. [\[CrossRef\]](#)
32. Derpich, M.S.; Silva, E.I.; Østergaard, J. Fundamental Inequalities and Identities Involving Mutual and Directed Informations in Closed-Loop Systems. *arXiv* **2013**, arXiv:1301.6427.
33. Shahsavari Baboukani, P.; Graversen, C.; Alickovic, E.; Østergaard, J. Estimating Conditional Transfer Entropy in Time Series Using Mutual Information and Nonlinear Prediction. *Entropy* **2020**, *22*, 1124. [\[CrossRef\]](#) [\[PubMed\]](#)
34. Barforooshan, M.; Derpich, M.S.; Stavrou, P.A.; Østergaard, J. The Effect of Time Delay on the Average Data Rate and Performance in Networked Control Systems. *IEEE Trans. Autom. Control* **2020**, *1*. [\[CrossRef\]](#)
35. Baboukani, P.S.; Graversen, C.; Østergaard, J. Estimation of Directed Dependencies in Time Series Using Conditional Mutual Information and Non-linear Prediction. In Proceedings of the 2020 28th European Signal Processing Conference (EUSIPCO), Amsterdam, The Netherlands, 18–21 January 2021; pp. 2388–2392. [\[CrossRef\]](#)
36. Yeh, J. *Real Analysis*, 3rd ed.; World Scientific: Singapore, 2014.

37. Gray, R.M. *Entropy and Information Theory*, 2nd ed.; Science+Business Media, Springer: New York, NY, USA, 2011.
38. Gray, R.M. *Probability, Random Processes and Ergodic Properties*, 2nd ed.; Springer: Berlin/Heidelberg, Germany, 2009.
39. Yeung, R.W. *A First Course in Information Theory*; Springer: Berlin/Heidelberg, Germany, 2002.
40. Goodwin, G.C.; Graebe, S.; Salgado, M.E. *Control System Design*; Prentice Hall: Hoboken, NJ, USA, 2001.
41. Sahai, A.; Mitter, S. The Necessity and Sufficiency of Anytime Capacity for Stabilization of a Linear System Over a Noisy Communication Link—Part I: Scalar Systems. *IEEE Trans. Inf. Theory* **2006**, *52*, 3369–3395. [[CrossRef](#)]
42. Khina, A.; Gårding, E.R.; Pettersson, G.M.; Kostina, V.; Hassibi, B. Control Over Gaussian Channels With and Without Source?Channel Separation. *IEEE Trans. Autom. Control* **2019**, *64*, 3690–3705. [[CrossRef](#)]
43. Khina, A.; Halbawi, W.; Hassibi, B. (Almost) practical tree codes. In Proceedings of the 2016 IEEE International Symposium on Information Theory, Barcelona, Spain, 10–15 July 2016; pp. 2404–2408. [[CrossRef](#)]