



**HAL**  
open science

## Quantitative Information Flow in Interactive Systems

Mário S. Alvim, Miguel E. Andrés, Catuscia Palamidessi

► **To cite this version:**

Mário S. Alvim, Miguel E. Andrés, Catuscia Palamidessi. Quantitative Information Flow in Interactive Systems. *Journal of Computer Security*, 2012, 20 (1), pp.3-50. inria-00637356

**HAL Id: inria-00637356**

**<https://inria.hal.science/inria-00637356>**

Submitted on 1 Nov 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Quantitative Information Flow in Interactive Systems

Mário S. Alvim<sup>1</sup>, Miguel E. Andrés<sup>2</sup>, and Catuscia Palamidessi<sup>1</sup>.

<sup>1</sup>INRIA and LIX, École Polytechnique Palaiseau, France.

<sup>2</sup>Institute for Computing and Information Sciences, The Netherlands.

**Abstract.** We consider the problem of defining the information leakage in interactive systems where secrets and observables can alternate during the computation. We show that the information-theoretic approach which interprets such systems as (simple) noisy channels is no longer valid. However, the principle can be recovered if we consider channels of a more complicated kind, that in Information Theory are known as channels with memory and feedback. We show that there is a complete correspondence between interactive systems and such kind of channels. Furthermore, we show that the capacity of the channels associated to such systems is a continuous function of a pseudometric based on the Kantorovich metric.

## 1 Introduction

Information leakage refers to the problem that arises when the observable behavior of a system reveals information that we would like to keep secret. This is also known as the problem of information flow from *high* variables to *low* variables. In recent years there has been a growing interest in quantitative approaches to this problem, because it is often desirable to quantify the partial knowledge of the secrets in terms of probability a distribution. Another reason is that the mechanisms to protect the information may use randomization to obfuscate the relation between the secrets and the observables.

Among the quantitative approaches, some of the most popular ones are based on Information Theory [5, 16, 4, 24, 6]. The idea is to interpret the system as an information-theoretic *channel*, where the secrets are the input and the observables are the output. The channel matrix consists of the conditional probabilities  $p(b | a)$ , defined as the measure of the executions producing the observable  $b$ , relatively to those which contain the secret  $a$ . The leakage is represented by the *mutual information*, and the worst-case leakage by the *capacity* of the channel.

In the above works, the secret value is assumed to be chosen at the beginning of the computation. We are interested in the more general scenario in which secrets can be chosen at any point. More precisely, we consider *interactive systems*, i.e. systems in which secrets and observables can alternate during the computation and influence each other. Examples of interactive systems include *auction protocols* like [31, 27, 25]. Some of these have become very popular thanks to their integration in Internet-based electronic commerce platforms [10, 11, 19]. Other examples of interactive programs include web servers, GUI applications, and command-line programs [3].

In this paper we investigate the applicability of the information-theoretic approach to interactive systems. In order to derive an information-theoretic channel, at a first look

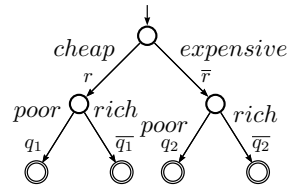
it would seem natural to define the matrix elements by using the definition of  $p(b|a)$  in terms of the joint and marginal probabilities  $p(a,b)$  and  $p(b)$ . Namely, the entry  $p(b|a)$  would be defined as the measure of the traces with (secret, observable)-projection  $(a,b)$ , divided by the measure of the traces with secret projection  $a$ . An approach of this kind was proposed in [9]. However, in the interactive case this construction does not really produce an information-theoretic channel. In fact, by definition a channel should be invariant with respect to the input distribution, and this is not the case here, as shown by the following example.

*Example 1.* Figure 1 represents a web-based interaction between one seller and two possible buyers, *rich* and *poor*. The seller offers two different products, *cheap* and *expensive*, with given probabilities. Once the product is offered, each buyer may try to buy it, with a certain probability. For simplicity we assume that the buyers offers are exclusive. We assume that the offers are observables, in the sense that they are made public on the website, while the identity of the buyer that actually buys the product should be kept secret from an external observer. The symbols  $r, q_1, q_2, \bar{r}, \bar{q}_1, \bar{q}_2$  represent probabilities, with the convention that  $\bar{r} = 1 - r$  (and the same for the pairs  $q_1, \bar{q}_1$  and  $q_2, \bar{q}_2$ ).

Following [9] we can compute the conditional probabilities as  $p(b|a) = \frac{p(a,b)}{p(a)}$ , thus obtaining the matrix in Table 1. However, the matrix is not invariant with respect to the input distribution. For instance, let us assume  $r = \bar{r} = \frac{1}{2}$ ,  $q_1 = \frac{2}{3}$ , and  $q_2 = \frac{2}{3}\rho$ , where  $\rho$  is a parameter. Therefore we have  $p(\text{poor}) = rq_1 + \bar{r}q_2 = \frac{1}{3}(1 + \rho)$  or, equivalently,  $\rho = 3 \cdot p(\text{poor}) - 1$ . Two different input distributions will determine different values of  $\rho$ , and therefore  $q_2$ . Hence also the channel matrices will be different, as the two examples in Table 2 show.

Consequently, when the secrets occur *after* the observables we cannot consider the conditional probabilities as representing a (classical) channel, and we cannot apply the standard information-theoretic concepts. In particular, we cannot adopt the (classical) capacity to represent the worst-case leakage, since the capacity is defined as the maximum information leakage using a fixed channel matrix over all possible input distributions. In other words, if we computed the (standard notion of) capacity using the matrix  $C$  obtained with a given input distribution  $D$  we could get a wrong result, because the capacity would in general correspond to the value of mutual information on a distribution  $D' \neq D$ . But since the matrix depends on the input distribution,  $D'$  would give a new matrix  $C' \neq C$ .

The first contribution of this paper is to consider an extension of the theory of channels which makes the information-theoretic approach applicable also the case of interactive systems. It turns out that a richer notion of channels, known in Information Theory as *channels with memory and feedback*, serves our purposes. The dependence of inputs on previous outputs corresponds to feedback, and the dependence of outputs



**Fig. 1.** An interactive syst.

	<i>cheap</i>	<i>expensive</i>
<i>poor</i>	$\frac{rq_1}{rq_1 + \bar{r}q_2}$	$\frac{\bar{r}q_2}{rq_1 + \bar{r}q_2}$
<i>rich</i>	$\frac{r\bar{q}_1}{rq_1 + \bar{r}q_2}$	$\frac{\bar{r}q_2}{rq_1 + \bar{r}q_2}$

**Table 1.** Channel matrix for Example 1

(a) $r = \frac{1}{2}, q_1 = \frac{2}{3}, \rho = \frac{1}{2}, q_2 = \frac{1}{3}$				(b) $r = \frac{1}{2}, q_1 = \frac{2}{3}, \rho = \frac{1}{4}, q_2 = \frac{1}{6}$			
	<i>cheap</i>	<i>expensive</i>	Input distr.		<i>cheap</i>	<i>expensive</i>	Input distr.
<i>poor</i>	$\frac{2}{3}$	$\frac{1}{3}$	$p(\textit{poor}) = \frac{1}{2}$	<i>poor</i>	$\frac{4}{5}$	$\frac{1}{5}$	$p(\textit{poor}) = \frac{5}{12}$
<i>rich</i>	$\frac{1}{3}$	$\frac{2}{3}$	$p(\textit{rich}) = \frac{1}{2}$	<i>rich</i>	$\frac{2}{7}$	$\frac{5}{7}$	$p(\textit{rich}) = \frac{7}{12}$

**Table 2.** Two different channel matrices induced by two different input distributions

on previous inputs and outputs corresponds to memory. Recent results in Information Theory [29] have shown that, in such channels, the transmission rate does not correspond to the maximum mutual information (the standard notion of capacity), but rather to the maximum normalized *directed information*, a concept introduced by Massey [17]. We propose to adopt this latter notion to represent leakage.

Our model of attacker is the interactive version of the attacker associated to Shannon entropy in the classification of Köpf and Basin [15]. We recall that in [15] an attacker is defined by the kind of questions that he can pose to an hypothetical oracle. In the case of Shannon entropy the questions are of the form “does  $s$  belong to  $S$ ?” where  $s$  is the secret that the attacker is trying to figure out, and  $S$  is a subset of the domain of secret values. The degree of invulnerability of the secret is the average number of questions that the attacker needs to ask in order to find out the exact value of the secret, under the best strategy (i.e. the best choice of the  $S$ ’s) for the given probability distribution on the secret values. It is easy to see that the invulnerability degree corresponds to the Shannon entropy of the secret. In the case of a standard single-use channel, the invulnerability degree of the secret *before* the attacker observes the output is the entropy of the input, determined by its a priori distribution. The invulnerability degree *after* the attacker observes the output is the conditional entropy of the input given the output, determined by its a posteriori distribution. The latter in general is lower than the first. The difference between these invulnerability degrees corresponds to the mutual information, and represents the leakage of the system.

In our interactive framework we consider the same scenario, but iterated. At each time step, we consider the input sequence so far; and the increase of its vulnerability caused by the observation of the new output is the contribution of the present step to the leakage. The sum of all these contributions represents the total leakage and, as we will see, corresponds to Massey’s directed information. We will come back to the model of attacker in Section 5, and discuss also a variant of this interpretation.

Gray investigated a concept similar to directed information, and he also conjectured the correspondence with the channel’s transmission rate [13]. His model is based on Millen’s synchronous state machines [20] and it is more general than ours, in that it admits observables and secrets at both ends of the channel. In other words, in addition to high inputs and low outputs, it considers also high inputs and low outputs. Gray derived his “quasi-directed-information” notion by extending Gallager’s formula for discrete finite state channels [12], and proposed it as definition of leakage, adducing as justification the above conjecture. However it is easy to see that the conjecture does not hold. We come back to this point in Section 5, after Definition 9.

A second contribution of our work is the proof that the channel capacity is a continuous function of a pseudometric on interactive systems based on the Kantorovich metric. The reason why we are interested in the continuity of the capacity is for computability purposes. Given a function  $f$  from a (pseudo)metric space  $X$  to a (pseudo)metric space  $Y$  the property of continuity for  $f$  means that, given a series of objects  $x_1, x_2, \dots \in X$  converging to  $x \in X$ , the series  $f(x_1), f(x_2), \dots \in Y$  converges to  $f(x) \in Y$ . Hence  $f(x)$  can be approximated by the objects  $f(x_1), f(x_2), \dots$ . The typical use of this property is in the case the trees are generated by programs containing loops. Generally the automaton expressing the semantics of the program can be seen as the (metric) limit of the sequence of trees generated by unfolding the loop at an increasingly deeper level. The continuity of the capacity means that we can approximate the real capacity by the capacities of these trees.

The continuity of the channel capacity was also proved in [9] for simple channels, but the proof does not adapt to the case of channels with memory and feedback and we had to devise a different technique. We illustrate this point by showing a counterexample (cfr. Example 6).

## 1.1 Plan of the paper

The paper is organized as follows. Section 2 reviews some important concepts from Probabilistic Automata and Information Theory. Section 3 reviews the notion of channel with memory and feedback that is the core of the model we propose. We discuss the concept of directed information and also the concept of capacity in the presence of feedback. Section 4 contains our main contribution. We explain how Interactive Information Hiding Systems (IIHSs) can be modeled using channels with memory and feedback. In particular we show that for any IIHS there is always a channel that simulates its probabilistic behavior. In Section 5 we discuss our notion of adversary and we define the quantification of information leakage as the channel's directed information from input to output, or as the directed capacity, depending on whether the input distribution is fixed or not. In Section 6 we show an example of our model applied to a protocol: the Cocaine Auction protocol. Section 7 proposes a pseudometric structure on IIHSs based on the Kantorovich metric. We also show that the capacity of the channels associated to interactive systems is a continuous function with respect to this pseudometric. In Sections 8 and 9 we review and discuss the main results of the paper and illustrate some future work.

A preliminary version of this paper appeared in the proceedings of CONCUR 2010 [1]. The additional material presented here consists in the proofs, the auxiliary Lemmata 2, 3, and 4, Propositions 2 and 3, more examples, and a more elaborate discussion about the model.

## 2 Preliminaries

In this section we briefly review some basic notions that we will need throughout the paper.

## 2.1 Probabilistic automata

A function  $\mu: \mathcal{S} \rightarrow [0, 1]$  is a *discrete probability distribution* on a countable set  $\mathcal{S}$  if  $\sum_{s \in \mathcal{S}} \mu(s) = 1$  and  $\mu(s) \geq 0$  for all  $s$ . The set of all discrete probability distributions on  $\mathcal{S}$  is  $\mathcal{D}(\mathcal{S})$ .

A *probabilistic automaton* [22] is a quadruple  $M = (\mathcal{S}, \mathcal{L}, \hat{s}, \vartheta)$  where  $\mathcal{S}$  is a countable set of *states*,  $\mathcal{L}$  a finite set of *labels* or *actions*,  $\hat{s}$  the *initial state*, and  $\vartheta$  a *transition function*  $\vartheta: \mathcal{S} \rightarrow \wp_f(\mathcal{D}(\mathcal{L} \times \mathcal{S}))$ . Here  $\wp_f(X)$  is the set of all finite subsets of  $X$ . If  $\vartheta(s) = \emptyset$  then  $s$  is a *terminal state*. We write  $s \rightarrow \mu$  for  $\mu \in \vartheta(s)$ ,  $s \in \mathcal{S}$ . Moreover, we write  $s \xrightarrow{\ell} r$  for  $s, r \in \mathcal{S}$  whenever  $s \rightarrow \mu$  and  $\mu(\ell, r) > 0$ . A *fully probabilistic automaton* is a probabilistic automaton satisfying  $|\vartheta(s)| \leq 1$  for all states. In such automata, when  $\vartheta(s) \neq \emptyset$ , we overload the notation and denote by  $\vartheta(s)$  the distribution outgoing from  $s$ .

A *path* in a probabilistic automaton is a sequence  $\sigma = s_0 \xrightarrow{\ell_1} s_1 \xrightarrow{\ell_2} \dots$  where  $s_i \in \mathcal{S}$ ,  $\ell_i \in \mathcal{L}$  and  $s_i \xrightarrow{\ell_{i+1}} s_{i+1}$ . A path can be *finite* in which case it ends with a state. A path is *complete* if it is either infinite, or finite ending in a terminal state. Given a finite path  $\sigma$ ,  $\text{last}(\sigma)$  denotes its last state. Let  $\text{Paths}_s(M)$  denote the set of all paths,  $\text{Paths}_s^*(M)$  the set of all finite paths, and  $\text{CPaths}_s(M)$  the set of all complete paths of an automaton  $M$ , starting from the state  $s$ . We will omit  $s$  if  $s = \hat{s}$ . Paths are ordered by the prefix relation, which we denote by  $\leq$ . The *trace* of a path is the sequence of actions in  $\mathcal{L}^* \cup \mathcal{L}^\infty$  obtained by removing the states, hence for the above  $\sigma$  we have  $\text{trace}(\sigma) = l_1 l_2 \dots$ . If  $\mathcal{L}' \subseteq \mathcal{L}$ , then  $\text{trace}_{\mathcal{L}'}(\sigma)$  is the projection of  $\text{trace}(\sigma)$  on the elements of  $\mathcal{L}'$ .

Let  $M = (\mathcal{S}, \mathcal{L}, \hat{s}, \vartheta)$  be a (fully) probabilistic automaton,  $s \in \mathcal{S}$  a state, and let  $\sigma \in \text{Paths}_s^*(M)$  be a finite path starting in  $s$ . The *cone* generated by  $\sigma$  is the set of complete paths  $\langle \sigma \rangle = \{\sigma' \in \text{CPaths}_s(M) \mid \sigma \leq \sigma'\}$ . Given a fully probabilistic automaton  $M = (\mathcal{S}, \mathcal{L}, \hat{s}, \vartheta)$  and a state  $s$ , we can calculate the *probability value*, denoted by  $\mathbf{P}_s(\sigma)$ , of any finite path  $\sigma$  starting in  $s$  as follows:  $\mathbf{P}_s(s) = 1$  and  $\mathbf{P}_s(\sigma \xrightarrow{\ell} s') = \mathbf{P}_s(\sigma) \mu(\ell, s')$ , where  $\text{last}(\sigma) \rightarrow \mu$ .

Let  $\Omega_s \triangleq \text{CPaths}_s(M)$  be the sample space, and let  $\mathcal{F}_s$  be the smallest  $\sigma$ -algebra generated by the cones. Then  $\mathbf{P}$  induces a unique *probability measure* on  $\mathcal{F}_s$  (which we will also denote by  $\mathbf{P}_s$ ) such that  $\mathbf{P}_s(\langle \sigma \rangle) = \mathbf{P}_s(\sigma)$  for every finite path  $\sigma$  starting in  $s$ . For  $s = \hat{s}$  we write  $\mathbf{P}$  instead of  $\mathbf{P}_{\hat{s}}$ .

Given a probability space  $(\Omega, \mathcal{F}, P)$  and two events  $A, B \in \mathcal{F}$  with  $P(B) > 0$ , the *conditional probability* of  $A$  given  $B$ ,  $P(A \mid B)$ , is defined as  $P(A \cap B)/P(B)$ .

## 2.2 Concepts from Information Theory

For more detailed information on this part we refer to [7]. Let  $A, B$  denote two random variables with corresponding probability distributions  $p_A(\cdot), p_B(\cdot)$ , respectively (we shall omit the subscripts when they are clear from the context). Let  $\mathcal{A} = \{a_1, \dots, a_n\}$ ,  $\mathcal{B} = \{b_1, \dots, b_m\}$  denote, respectively, the sets of possible values for  $A$  and for  $B$ .

The *entropy* of  $A$  is defined as  $H(A) = -\sum_{a \in \mathcal{A}} p(a) \log p(a)$  and it measures the uncertainty of  $A$ . It takes its minimum value  $H(A) = 0$  when  $p_A(\cdot)$  is a point mass (also called delta of Dirac). The maximum value  $H(A) = \log |\mathcal{A}|$  is obtained when  $p_A(\cdot)$  is the uniform distribution. Usually the base of the logarithm is set to

be 2 and the entropy is measured in *bits*. The *conditional entropy* of  $A$  given  $B$  is  $H(A|B) = -\sum_B p(b) \sum_A p(a|b) \log p(a|b)$ , and it measures the uncertainty of  $A$  when  $B$  is known. It is well-known that  $0 \leq H(A|B) \leq H(A)$ . The minimum value, 0, is obtained when  $A$  is completely determined by  $B$ . The maximum value  $H(A)$  is obtained when  $A$  and  $B$  are independent. The *mutual information* between  $A$  and  $B$  is defined as  $I(A; B) = H(A) - H(A|B)$ , and it measures the amount of information about  $A$  that we gain by observing  $B$ . It can be shown that  $I(A; B) = I(B; A)$  and  $0 \leq I(A; B) \leq H(A)$ . If  $C$  is a third random variable, the *conditional mutual information* between  $A$  and  $B$  given  $C$  is defined as  $I(A; B|C) = H(A|C) - H(A|B, C)$ .

The (conditional) entropy and mutual information respect the *chain rules*. Namely, given the random variables  $A_1, A_2, \dots, A_k, B$  and  $C$ , we have:

$$H(A_1, A_2, \dots, A_k|C) = \sum_{i=1}^k H(A_i|A_1, \dots, A_{i-1}, C) \quad (1)$$

$$I(A_1, A_2, \dots, A_k; B|C) = \sum_{i=1}^k I(A_i; B|A_1, \dots, A_{i-1}, C) \quad (2)$$

a family  $\rho = \{p_v(\cdot)\}_v$  of probability measures parametrized on  $v$  is called a *stochastic kernel*<sup>1</sup>.

A (*discrete memoryless*) *channel* is a tuple  $(\mathcal{A}, \mathcal{B}, p(\cdot|\cdot))$ , where  $\mathcal{A}, \mathcal{B}$  are the sets of input and output symbols, respectively, and  $p(b|a)$  is the probability of observing the output symbol  $b$  when the input symbol is  $a$ . These conditional probabilities constitute the *channel matrix*. An input distribution  $p_A(\cdot)$  over  $\mathcal{A}$  together with the channel determine the joint distribution  $p(a, b) = p(a|b) \cdot p(a)$  and consequently  $I(A; B)$ . The maximum  $I(A; B)$  over all possible input distributions is the channel's *capacity*.

### 3 Discrete channels with memory and feedback

In this section we present the notion of channel with memory and feedback. We assume a scenario in which the channel is used repeatedly, in a finite temporal sequence of steps  $1, \dots, T$ . Intuitively, memory means that the output at time  $t$  depends on the input and output histories, i.e. on the inputs till time  $t$ , and on the output till time  $t - 1$ . Feedback means that the input at time  $t$  depends on the outputs till time  $t - 1$ .

We adopt the following notation, which appears to be standard in the literature of channels with memory and feedback.

**Convention 1.** Given a set of symbols (alphabet)  $\mathcal{A} = \{a_1, \dots, a_n\}$ , we use a Greek letter  $(\alpha, \beta, \dots)$  to denote a sequence of symbols ordered in time. Given a sequence  $\alpha = a_{i_1} a_{i_2} \dots a_{i_m}$ , the notation  $\alpha_t$  represents the symbol at time  $t$ , i.e.  $a_{i_t}$ , while  $\alpha^t$  represents the sequence  $\alpha_{i_1} \alpha_{i_2} \dots \alpha_{i_t}$ . For instance, in the sequence  $\alpha = a_3 a_7 a_5$ , we have  $\alpha_2 = a_7$  and  $\alpha^2 = a_3 a_7$ . Analogously, if  $X$  is a random variable, then  $X^t$  denotes the sequence of  $t$  consecutive instances  $X_1, \dots, X_t$  of  $X$ .

<sup>1</sup> The general definition of stochastic kernel is more complicated (cfr. [29]), but it reduces to this one in the case of discrete channels, which is what we use in this paper.

We now define formally the concepts of memory and feedback. Consider a channel from input  $A$  to output  $B$ . The channel behavior after  $T$  uses can be fully described by the joint distribution of  $A^T \times B^T$ , namely by the probabilities  $p(\alpha^T, \beta^T)$ . Using the chain rule, we can decompose these probabilities as follows:

$$p(\alpha^T, \beta^T) = \prod_{t=1}^T p(\alpha_t | \alpha^{t-1}, \beta^{t-1}) p(\beta_t | \alpha^t, \beta^{t-1}) \quad (3)$$

**Definition 1.** We say that the channel has feedback if, in general,  $p(\alpha_t | \alpha^{t-1}, \beta^{t-1}) \neq p(\alpha_t | \alpha^{t-1})$ , i.e. the probability of  $\alpha_t$  depends, besides  $\alpha^{t-1}$ , also on  $\beta^{t-1}$ . Analogously, we say that the channel has memory if, in general,  $p(\beta_t | \alpha^t, \beta^{t-1}) \neq p(\beta_t | \alpha_t)$ , i.e. the probability of  $\beta_t$  depends on  $\alpha^t$  and  $\beta^{t-1}$ .

Note that in the opposite case, i.e. when  $p(\alpha_t | \alpha^{t-1}, \beta^{t-1})$  coincides with  $p(\alpha_t | \alpha^{t-1})$  and  $p(\beta_t | \alpha^t, \beta^{t-1})$  coincides with  $p(\beta_t | \alpha_t)$ , then we have a classical channel (memoryless, and without feedback), in which each use is independent from the previous ones. The only possible dependency on the history is the one of  $\alpha_t$  on  $\alpha^{t-1}$ . This is because  $A_1, \dots, A_T$  are in general correlated, due to the fact that they are produced by an encoding function. Note that in absence of memory and feedback (3) reduces to  $p(\alpha^T, \beta^T) = \prod_{t=1}^T p(\alpha_t, \beta_t) p(\beta_t | \alpha_t)$ , which is the standard formula for a classical channel after  $T$  uses.

The above is a very abstract description of a channel with memory and feedback. We now discuss a more concrete notion following the presentation of [29]. Such a channel, represented in Figure 2, consists of a sequence of components formally defined as a family of stochastic kernels  $\{p(\cdot | \alpha^t, \beta^{t-1})\}_{t=1}^T$  over  $\mathcal{B}$ . The probabilities  $p(\beta_t | \alpha^t, \beta^{t-1})$  represent the channel *innermost behavior* at time  $t$ ,  $1 \leq t \leq T$ : the internal channel takes the input  $\alpha_t$  and, depending the history of inputs and outputs so far, it produces an output symbol  $\beta_t$ . The output is then fed back to the encoder with delay one. On the output side, at time  $t$  the encoder takes the message and the past output symbols  $\beta^{t-1}$  and produces a channel input symbol  $\alpha_t$  according to the code function  $\varphi_t$ . At final time  $T$  the decoder takes all the channel outputs  $\beta^T$  and produces the decoded message  $\hat{W}$ . The order is the following:

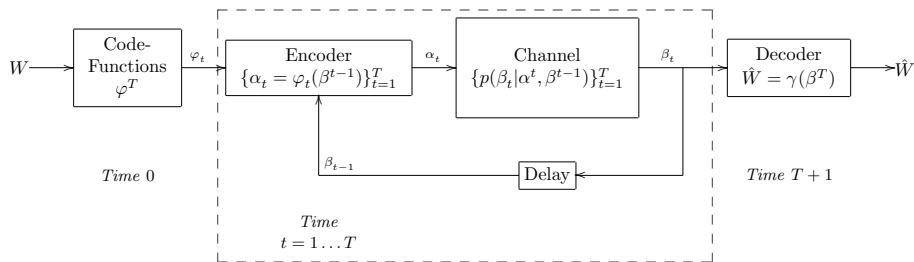
$$\text{Message } W, \quad \alpha_1, \beta_1, \quad \alpha_2, \beta_2, \quad \dots, \quad \alpha_T, \beta_T, \quad \text{Decoded Message } \hat{W} \quad (4)$$

Let us now explain the concept of code function. Intuitively, a code function is a strategy to encode the message into a suitable representation to be transmitted through the channel. There is a code function for each possible message, and the function is fixed at the very beginning of the transmission (time  $t = 0$ ). However, the encoding can use the information provided via feedback, so each component  $\varphi_t$  ( $1 \leq t \leq T$ ) of the code function takes as parameter the history of feedback  $\beta^{t-1}$  to generate the next input symbol  $\alpha_t$ .

Formally, let  $\mathcal{F}_t$  be the set of all measurable maps  $\varphi_t : \mathcal{B}^{t-1} \rightarrow \mathcal{A}$  endowed with a probability distribution, and let  $F_t$  be the corresponding random variable. Let  $\mathcal{F}^T, F^T$  denote the Cartesian product on the domain and the random variable, respectively. A *channel code function* is an element  $\varphi^T = (\varphi_1, \dots, \varphi_T) \in \mathcal{F}^T$ .

MS: Discuss about this notation: is it correctly denoting a sequence of matrices, one for each time  $t$ ? I'd say it stands for  $\{p_{A_t | A^{t-1}, B^{t-1}}(\cdot | \cdot)\}_{t=1}^T$ .





**Fig. 2.** Model for discrete channel with memory and feedback

Note that, by the chain rule,  $p(\varphi^T) = \prod_{t=1}^T p(\varphi_t | \varphi^{t-1})$ . Hence the distribution on  $\mathcal{F}^T$  is uniquely determined by a sequence  $\{p(\varphi_t | \varphi^{t-1})\}_{t=1}^T$ . We will use the notation  $\varphi^t(\beta^{t-1})$  to represent the  $\mathcal{A}$ -valued  $t$ -tuple  $(\varphi_1, \varphi_2(\beta^1), \dots, \varphi_t(\beta^{t-1}))$ .

In Information Theory this kind of channel is used to encode and transmit messages. If  $\mathcal{W}$  is a set of messages of cardinality  $M$  with typical element  $w$ , endowed with a probability distribution, a *channel code* is a set of  $M$  channel code functions  $\varphi^T[w]$ , interpreted as follows: for message  $w$ , if at time  $t$  the channel feedback is  $\beta^{t-1}$ , then the channel encoder outputs  $\varphi_t[w](\beta^{t-1})$ . A *channel decoder* is a map from  $\mathcal{B}^T$  to  $\mathcal{W}$  which attempts to reconstruct the input message after observing all the output history  $\beta^T$  from the channel.

### 3.1 The power of feedback

The original purpose of *communication channels* models is to represent data transmission from a source to a receiver. Shannon's Channel Coding Theorem states for every channel there is an encoding scheme that allows a transmission rate arbitrary close to the channel capacity with a negligible probability of error (if the number of uses of the channel is large enough). Shannon did not explain however how to determine such an encoding, and a general way to generate an optimal encoding scheme has not been found yet. The use of feedback, fortunately, can simplify the design of the encoder and of the decoder. The following example illustrates the idea.

*Example 2.* Consider a discrete memoryless binary channel  $\{\mathcal{A}, \mathcal{B}, p(\cdot|\cdot)\}$  with  $\mathcal{A} = \{0, 1\}$ ,  $\mathcal{B} = \{0, 1, e\}$  and the channel matrix of Table 3. This kind of channel is called *erasure channel* because it can lose (or *erase*) bits during the transmission with a certain probability. Namely, any bit has 0.8 probability of being correctly transmitted, and 0.2 probability of being lost. On the output side the encoder is able to detect whether the bit was erased (by receiving an  $e$  symbol), but it cannot tell which was the actual value of the original bit. The Channel Coding Theorem guarantees that the maximum information transmission rate in this channel is (2 to the power of) the channel capacity, i.e 0.8 bits per use of the channel.

	0	1	e
0	0.8	0	0.2
1	0	0.8	0.2

**Table 3.** Channel matrix for binary erasure channel

Following simple principles described in [7], an encoding that achieves the capacity can be easily obtained if the channel can be used with feedback. The idea is an adaptation of the stop-and-wait protocol [26, 28]. Suppose that every bit received on the output end of the channel is fed back noiselessly to the source with delay 1. Define the encoding as follows: for each bit transmitted, the encoder checks via feedback if the bit was erased. If not, the encoder moves on to transmit the next of the message. If yes, the encoder transmits the same bit again.

It is easy to see that with this encoding scheme the transmission rate is 0.8 bit per usage of the channel, since in 80% of the cases the bit is transmitted properly, and in 20% it is lost and a retransmission is needed.

In the appendix (Section 9) we come back to this example to illustrate more in detail the design and the function of the encoder and decoder.

Note that the channel capacity in the above example does not increase with the addition of feedback (it is 0.8 bit per usage of the channel with or without feedback). This is because the channel is memoryless: *feedback does not increase the capacity of discrete memoryless channels* [7]. In general however, feedback *does* increase the capacity.

### 3.2 Directed information and capacity of channels with feedback

In classical Information Theory, the channel capacity, which is related to the channel's transmission rate by Shannon's Channel Coding Theorem, can be obtained as the supremum of the mutual information over all possible input distributions. In the presence of feedback, however, this correspondence does not longer hold. More specifically, mutual information no longer represents the information flow from  $A^T$  to  $B^T$ . Intuitively, this is due to the fact that mutual information expresses correlation, and therefore it is increased by feedback. However, feedback, i.e the way the output influences the next input, is not part of the information to be transmitted. If we want to maintain the correspondence between the transmission rate and capacity, we need to replace the mutual information with *directed information* [17].

**Definition 2.** *In a channel with feedback, the directed information from input  $A^T$  to output  $B^T$  is defined as  $I(A^T \rightarrow B^T) = \sum_{t=1}^T I(A^t; B_t | B^{t-1})$ . In the other direction, the directed information from  $B^T$  to  $A^T$  is defined as:  $I(B^T \rightarrow A^T) = \sum_{t=1}^T I(A_t; B^{t-1} | A^{t-1})$ .*

In Section 5 we shall discuss relation between directed information and mutual information, as well as the correspondence with information leakage. For the moment, we only present the extension of the concept of capacity.

Let  $\mathcal{D}_T = \{\{p(\alpha_t | \alpha^{t-1}, \beta^{t-1})\}_{t=1}^T\}$  be the set of all input distributions in presence of feedback. For finite  $T^2$ , the capacity of a channel with memory and feedback is:

$$C_T = \sup_{\mathcal{D}_T} \frac{1}{T} I(A^T \rightarrow B^T) \quad (5)$$

<sup>2</sup> For infinite  $T$ , see the definition in [29]. However this definition is not used in this paper.

## 4 Interactive systems as channels with memory and feedback

Interactive Information Hiding Systems (IIHS) [2] are a variant of probabilistic automata in which we separate actions into secrets (inputs) and observables (outputs). “Interactive” means that secrets and observables can interleave and influence each other. In this paper we consider only IIHSs of two particular kinds: the *fully probabilistic* IIHSs, where there is no nondeterminism, and the *secret-nondeterministic* (or input-nondeterministic) IIHSs, where each secret choice is fully nondeterministic.

In this section we formalize the notion of IIHS and we show how to associate to an IIHS a channel with memory and feedback.

**Definition 3.** An IIHS is a triple  $\mathcal{J} = (M, \mathcal{A}, \mathcal{B})$ , where  $\mathcal{A}$  and  $\mathcal{B}$  are disjoint sets of secrets and observables respectively,  $M$  is a probabilistic automaton  $(\mathcal{S}, \mathcal{L}, \hat{s}, \vartheta)$  with  $\mathcal{L} = \mathcal{A} \cup \mathcal{B}$ , and, for each  $s \in \mathcal{S}$ :

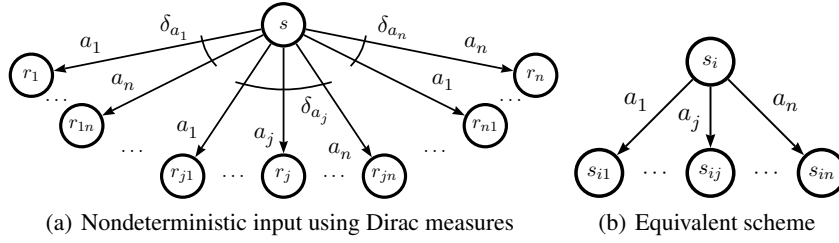
1. either  $\vartheta(s) \subseteq \mathcal{D}(\mathcal{A} \times \mathcal{S})$  ( $s$  is a secret state) or  $\vartheta(s) \subseteq \mathcal{D}(\mathcal{B} \times \mathcal{S})$  ( $s$  is an observable state)
2. if  $s \xrightarrow{\ell} r$  then: if  $s$  is secret then  $r$  is observable, and if  $s$  is observable then  $r$  is secret
3. if  $\vartheta(s) \subseteq \mathcal{D}(\mathcal{B} \times \mathcal{S})$  then  $|\vartheta(s)| \leq 1$
4. if  $\vartheta(s) \subseteq \mathcal{D}(\mathcal{A} \times \mathcal{S})$  then either
  - $|\vartheta(s)| \leq 1$  (fully probabilistic IIHS) or
  - there exist  $a'_i s$  and  $s'_i s$  ( $i = 1, \dots, n$ ) such that  $\vartheta(s) = \{\delta(a_i, s_i)\}_{i=1}^n$ , where  $\delta(a_i, s_i)$  is the Dirac measure (secret-nondeterministic IIHS).

In the above definition, Conditions 1 and 2 imply that the IIHS is alternating between secrets and observables. Once unfolded, all the transitions between nodes at two consecutive depths have either secret actions only, or observable actions only. Moreover, the occurrences of secret and observable actions alternate. We also assume without loss of generality that the first level contains secret actions. We call *normalized* the automata that satisfy these properties. We note that in our context this is not really a restriction, because given a IIHS which is not normalized, it is always possible to transform it into a normalized IIHS which is equivalent to the former one up to a given execution level. The reader can find in the appendix (Section 9) the formal definition of the transformation.

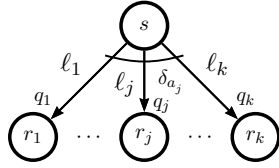
Note that Condition 3 means that all observable transitions are fully probabilistic. Condition 4 means that all secret transitions are either fully probabilistic or fully nondeterministic. The latter case is justified by the fact that secret-nondeterministic IIHS the secret transition scheme, represented in Figure 3(a), is equivalent to the one of Figure 3(b), where every possible action is allowed to lead to exactly one state.

Note that we do not consider here internal nondeterminism such as, for instance, that arising from interleaving of concurrent processes. This means that we make a rather restricted use of the notion of Probabilistic Automaton, but this is enough for our purposes. The presence of nondeterminism gives rise to a new set of problems (see for example [4]) which are orthogonal to those considered in this paper.

We show now that the secret and observable traces determine the states, hence they are enough to retrieve the path. We first need the following auxiliary lemma:



**Fig. 3.** Scheme of secret transitions for secret-nondeterministic IIHSs



**Fig. 4.** Typical transition in a fully probabilistic IIHS

**Lemma 1.** *Given an IIHS, for every  $s$  and  $\ell$  there exists a unique  $r$  such that  $s \xrightarrow{\ell} r$ .*

*Proof.* First we observe that, for both kinds of IIHSs, if  $s$  is a secret state then the property follows immediately from the definition of IIHSs, Condition 1 (see Figure 4). The case of secret state for a fully probabilistic IIHS is analogous. For secret-nondeterministic IIHSs, from a secret state there may be several possible outgoing probability distributions that can be nondeterministically chosen. However, every possible distribution is a Dirac measure of a different secret symbol. This means that there is at most one way of performing a transition under some specific action.  $\square$

**Proposition 1.** *Given an IIHS, consider two paths  $\sigma$  and  $\sigma'$ . If  $\text{trace}_{\mathcal{A}}(\sigma) = \text{trace}_{\mathcal{A}}(\sigma')$  and  $\text{trace}_{\mathcal{B}}(\sigma) = \text{trace}_{\mathcal{B}}(\sigma')$ , then  $\sigma = \sigma'$ .*

*Proof.* By induction on the length of the traces. The initial state of the automaton is uniquely determined by the empty (secret and observable) traces. Assume now we are in a state  $s$  uniquely determined by secret and observable traces  $\alpha$  and  $\beta$ , respectively. If  $s$  makes a secret transition  $s \xrightarrow{a} s'$ , then by Lemma 1 there is only one state  $s'$  reachable from  $s$  via an  $a$ -transition, and therefore  $s'$  is uniquely determined by the secret trace  $\alpha' = \alpha a$  and the observable trace  $\beta$ . The case in which  $s$  makes an observable transition is similar.  $\square$

#### 4.1 Construction of the channel associated to an IIHS

We now show how to associate a channel to an IIHS.

In an interactive systems secrets and observables may interleave and influence each other. Considering a channel with memory and feedback is a way to capture this rich behavior. Secrets have a causal influence on observables via the channel, and, in the presence of interactivity, observables have a causal influence on secrets via feedback. This alternating mutual influence between secrets and observables can be modeled by repeated uses of the channels. Each time the channel is used it represents a different

state of the computation, and the conditional probabilities of observables on secrets can depend on this state. The addition of memory to the model allows expressing the dependency of the channel matrix on such a state.

We will see that a secret-nondeterministic IIHS determines a channel as specified by its stochastic kernel, while a fully probabilistic IIHS determines, additionally, the input distribution.

In Section 6 we will give an extensive and detailed example of how to make such a construction for a real security protocol.

Given a path  $\sigma$  of length  $2t - 1$ , we will denote  $\text{trace}_{\mathcal{A}}(\sigma)$  by  $\alpha^t$ , and  $\text{trace}_{\mathcal{B}}(\sigma)$  by  $\beta^{t-1}$ .

**Definition 4.** *Let  $\mathcal{J}$  be an IIHS. For each  $t$ , the channel's stochastic kernel corresponding to  $\mathcal{J}$  is defined as  $p(\beta_t | \alpha^t, \beta^{t-1}) = \vartheta(s)(\beta_t, s')$ , where  $s$  is the state reached from the root via the path  $\sigma$  whose secret and observable trace are  $\alpha^t$  and  $\beta^{t-1}$  respectively.*

Note that  $s$  and  $s'$  in the previous definition are well defined: by Proposition 1,  $s$  is unique, and since the choice of  $\beta_t$  is fully probabilistic,  $s'$  is also unique.

The following example illustrates how to apply Definition 4, with the help of Proposition 1, to build the channel matrix of a simple example.

*Example 3.* Let us consider an extended version of the website interactive system of Figure 1. We maintain the general definition of the system, i.e, there are two possible buyers (*rich* and *poor* represented by *rc.* and *pr.*, respectively) and two possible products (*cheap* and *expensive*, represented by *chp.* and *exp.*, respectively). We still assume that offers are observable, since they are visible to everyone on the website, but the identity of buyers should be kept secret. We consider two consecutive rounds of offers and buys, which implies that, after normalization,  $T = 3$ . Figure 5 shows an automaton for this example in normalized form. Transitions with null probability are omitted, and the symbol  $a_*$  is used as a place holder to achieve the normalized IIHS (see Appendix).

To construct the stochastic kernels  $\{p(\beta_t | \alpha^t, \beta^{t-1})\}_{t=1}^T$ , we need to determine the conditional probability of an observable at time  $t$  given the history up to time  $t$ .

Let us take the case  $t = 2$  and compute the conditional probability of the observable  $\beta_2 = \textit{cheap}$  given that the history of secrets until time  $t = 2$  is  $\alpha^2 = a_*, \textit{poor}$  and the history of observables is  $\beta^1 = \textit{expensive}$ . Applying Definition 4, we see that  $p(\beta_2 = \textit{cheap} | \alpha^2 = a_*, \textit{poor}, \beta^1 = \textit{expensive}) = \vartheta(s)(\textit{cheap}, s')$ . By Proposition 1, the traces  $\alpha^2 = a_*, \textit{poor}, \beta^1 = \textit{expensive}$  determine a unique state  $s$  in the automaton, namely, the state  $s = 5$ . Moreover, from the state 5 a unique transition labelled with the action *cheap* is possible, leading to the state  $s' = 11$ . Therefore, we can conclude that  $p(\beta_2 = \textit{cheap} | \alpha^2 = a_*, \textit{poor}, \beta^1 = \textit{expensive}) = \vartheta(s = 5)(\textit{cheap}, s' = 11) = p_{23}$ .

Similarly, with  $t = 1$  and history  $\alpha^1 = a_*, \beta^0 = \epsilon$ , the observable symbol  $\beta_1 = \textit{expensive}$  can be observed with probability  $p(\beta_1 = \textit{expensive} | \alpha^1 = a_*, \beta^0 = \epsilon) = \vartheta(s = 0)(\textit{cheap}, s' = 2) = \bar{p}_1$ .

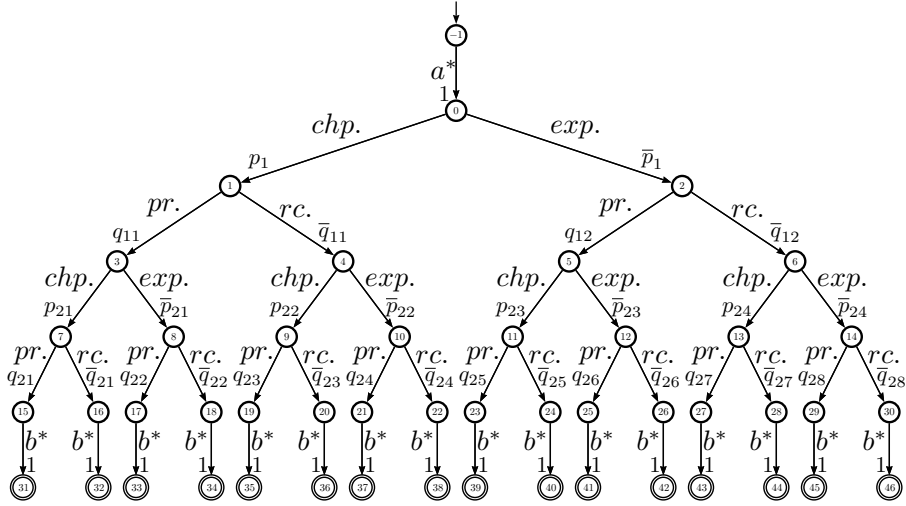
If  $\mathcal{J}$  is fully probabilistic, then it determines also the input distribution and the dependency of  $\alpha_t$  on  $\beta^{t-1}$  (feedback) and on  $\alpha^{t-1}$ .

**Definition 5.** Let  $\mathcal{J}$  be an IIHS. If  $\mathcal{J}$  is fully probabilistic, the associated channel has a conditional input distribution for each  $t$  defined as  $p(\alpha_t|\alpha^{t-1}, \beta^{t-1}) = \vartheta(s)(\alpha_t, s')$ , where  $s$  is the state reached from the root via the path  $\sigma$  whose secret and observable traces are  $\alpha^{t-1}$  and  $\beta^{t-1}$  respectively.

*Example 4.* Since the system of Example 3 is fully probabilistic, we can calculate the values of the conditional probabilities  $\{p(\alpha_t|\alpha^{t-1}, \beta^{t-1})\}_{t=1}^T$ .

Let us take, for instance, the case where  $t = 2$  and compute the conditional probability of secret  $\alpha_2 = \text{poor}$  given that the history of secrets until time  $t = 2$  is  $\alpha^1 = a_*$  and the history of observables is  $\beta^1 = \text{expensive}$ . Applying Definition 5, we see that  $p(\alpha_2 = \text{poor}|\alpha_1 = a_*, \beta^1 = \text{expensive}) = \vartheta(s)(\text{poor}, s')$ . By Proposition 1, the traces  $\alpha^1 = a_*, \beta^1 = \text{expensive}$  determine a unique state  $s$  in the automaton, namely, the state  $s = 2$ . Moreover, from the state 2 a unique transition labelled with the action  $\text{poor}$  is possible, leading to the state  $s' = 5$ . Therefore, we can conclude that  $p(\alpha_2 = \text{poor}|\alpha_1 = a_*, \beta^1 = \text{expensive}) = \vartheta(s = 2)(\text{poor}, s' = 5) = q_{12}$ .

Similarly, with  $t = 3$  and history  $\alpha^2 = a_*, \text{rich}, \beta^2 = \text{cheap}, \text{expensive}$ , the secret symbol  $\alpha_3 = \text{rich}$  can be observed with probability  $p(\alpha_3 = \text{rich}|\alpha^2 = a_*, \text{rich}, \beta^2 = \text{cheap}, \text{expensive}) = \vartheta(s = 10)(\text{cheap}, s' = 22) = \bar{q}_{24}$ .



**Fig. 5.** The normalized IIHS for the extended website example

## 4.2 Lifting the channel inputs to reaction functions

Definitions 4 and 5 show how to obtain the the joint probabilities  $p(\alpha^t, \beta^t)$  for a fully probabilistic IIHS. We still need to show in what sense this joint probability distribution defines an information-theoretic channel.

The  $\{p(\beta_t|\alpha^t, \beta^{t-1})\}_{t=1}^T$  determined by the IIHS correspond to a channel's stochastic kernel. The problem resides in the conditional probabilities  $\{p(\alpha_t|\alpha^{t-1}, \beta^{t-1})\}_{t=1}^T$ . In an information-theoretic channel, the value of  $\alpha_t$  is determined in the encoder by a deterministic function  $\varphi_t(\beta^{t-1})$ . Therefore, inside the encoder there is no possibility for a probabilistic description of  $\alpha_t$ . The solution is to externalize this probabilistic behavior to the code functions.

As shown in [29], the original channel with feedback from input symbols  $\mathcal{A}^T$  to output symbols  $\mathcal{B}^T$  can be lifted to an equivalent channel without feedback from code functions  $\mathcal{F}^T$  to output symbols  $\mathcal{B}^T$ . This transformation also allows us to calculate the channel capacity. Let  $\{p(\varphi_t|\varphi^{t-1})\}_{t=1}^T$  be a sequence of code function stochastic kernels and let  $\{p(\beta_t|\alpha^t, \beta^{t-1})\}_{t=1}^T$  be a channel with memory and feedback. The channel from  $F^T$  to  $B^T$  is constructed using a joint measure  $Q(\varphi^T, \alpha^T, \beta^T)$  that respects the following constraints:

**Definition 6.** A measure  $Q(\varphi^T, \alpha^T, \beta^T)$  is said to be consistent with respect to the code function stochastic kernels  $\{p(\varphi_t|\varphi^{t-1})\}_{t=1}^T$  and the channel  $\{p(\beta_t|\alpha^t, \beta^{t-1})\}_{t=1}^T$  if, for each  $t$ :

1. There is no feedback to the code functions:  $Q(\varphi_t|\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) = p(\varphi_t|\varphi^{t-1})$ .
2. The input is a function of the past outputs:  $Q(\alpha_t|\varphi^t, \alpha^{t-1}, \beta^{t-1}) = \delta_{\{\varphi_t(\beta^{t-1})\}}(\alpha_t)$  where  $\delta$  is the Dirac measure.
3. The properties of the underlying channel are preserved:

$$Q(\beta_t|F^t = \varphi^t, A^t = \alpha^t, B^{t-1} = \beta^{t-1}) = p(\beta_t|\alpha^t, \beta^{t-1})$$

The following result states that there is only one consistent measure  $Q(\varphi^T, \alpha^T, \beta^T)$ :

**Theorem 2 ([29]).** Given  $\{p(\varphi_t|\varphi^{t-1})\}_{t=1}^T$  and a channel  $\{p(\beta_t|\alpha^t, \beta^{t-1})\}_{t=1}^T$ , there exists only one consistent measure  $Q(\varphi^T, \alpha^T, \beta^T)$ . Furthermore the channel from  $\mathcal{F}^T$  to  $\mathcal{B}^T$  is given by:

$$Q(\beta_t|\varphi^t, \beta^{t-1}) = p(\beta_t|\varphi^t(\beta^{t-1}), \beta^{t-1}) \quad (6)$$

Since in our setting the concept of encoder makes little sense as there is no information to encode, we externalize the probabilistic behavior of  $\alpha_t$  as follows. Code functions become a single set of reaction functions  $\{\varphi_t\}_{t=1}^T$  with  $\beta^{t-1}$  as parameter (the message  $w$  does not play a role any more). Reaction functions can be seen as a model of how the environment reacts to given system outputs, producing new system inputs (they do not play a role of encoding a message). These reaction functions are endowed with a probability distribution that generates the probabilistic behavior of the values of  $\alpha_t$ .

**Definition 7.** A reactor is a distribution on reaction functions, i.e., a stochastic kernel  $\{p(\varphi_t|\varphi^{t-1})\}_{t=1}^T$ . A reactor  $R$  is consistent with a fully probabilistic IIHS  $\mathcal{I}$  if it induces the compatible distribution  $Q(\varphi^T, \alpha^T, \beta^T)$  such that, for every  $1 \leq t \leq T$ ,  $Q(\alpha_t|\alpha^{t-1}, \beta^{t-1}) = p(\alpha_t|\alpha^{t-1}, \beta^{t-1})$ , where the latter is the probability distribution induced by  $\mathcal{J}$ .

The main result of this section states that for any fully probabilistic IIHS there is a reactor that generates the probabilistic behavior of the IIHS.

**Lemma 2.** Let  $\mathcal{X}, \mathcal{Y}$  be non-empty finite sets, and let  $\tilde{x} \in \mathcal{X}, \tilde{y} \in \mathcal{Y}$ . Let  $p : \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$  be a function such that, for every  $x \in \mathcal{X}$ , we have:  $\sum_{y \in \mathcal{Y}} p(x, y) = 1$ . Then:

$$\sum_{\substack{f \in \mathcal{X} \rightarrow \mathcal{Y} \\ f(\tilde{x}) = \tilde{y}}} \prod_{x \in \mathcal{X}} p(x, f(x)) = p(\tilde{x}, \tilde{y})$$

*Proof.* By induction on the number of elements of  $\mathcal{X}$ .

**Base case:**  $\mathcal{X} = \{\tilde{x}\}$ . In this case:

$$\sum_{\substack{f \in \mathcal{X} \rightarrow \mathcal{Y} \\ f(\tilde{x}) = \tilde{y}}} \prod_{x \in \mathcal{X}} p(x, f(x)) = p(\tilde{x}, f(\tilde{x})) = p(\tilde{x}, \tilde{y})$$

**Inductive case:** Let  $\mathcal{X} = \mathcal{X}' \cup \{\hat{x}\}$ , with  $\tilde{x} \in \mathcal{X}'$  and  $\hat{x} \notin \mathcal{X}'$ . Then:

$$\begin{aligned} & \sum_{\substack{f \in \mathcal{X}' \cup \{\hat{x}\} \rightarrow \mathcal{Y} \\ f(\tilde{x}) = \tilde{y}}} \prod_{x \in \mathcal{X}' \cup \{\hat{x}\}} p(x, f(x)) \\ &= \text{(by distributivity)} \\ & \left( \sum_{\substack{f \in \mathcal{X}' \rightarrow \mathcal{Y} \\ f(\tilde{x}) = \tilde{y}}} \prod_{x \in \mathcal{X}'} p(x, f(x)) \right) \cdot \sum_{g \in \{\hat{x}\} \rightarrow \mathcal{Y}} p(\hat{x}, g(\hat{x})) \\ &= \text{(by the assumption)} \\ & \sum_{\substack{f \in \mathcal{X}' \rightarrow \mathcal{Y} \\ f(\tilde{x}) = \tilde{y}}} \prod_{x \in \mathcal{X}'} p(x, f(x)) \\ &= \text{(by the induction hypothesis)} \\ & p(\tilde{x}, \tilde{y}) \end{aligned}$$

□

**Theorem 3.** Let  $\mathcal{J}$  be a fully probabilistic IHS inducing the joint probability distribution  $p(\alpha^t, \beta^t)$ ,  $1 \leq t \leq T$ , on secret and observable traces. It is always possible to construct a channel with memory and feedback, and an associated probability distribution  $Q(\varphi^T, \alpha^T, \beta^T)$ , which corresponds to  $\mathcal{J}$  in the sense that, for every  $1 \leq t \leq T$ ,  $\alpha^t, \beta^t$ , the equality  $Q(\alpha^t, \beta^t) = p(\alpha^t, \beta^t)$  holds.

*Proof.* First of all we note that, by probability laws,  $Q(\alpha^t, \beta^t) = \sum_{\varphi^t} Q(\varphi^t, \alpha^t, \beta^t)$ . So we need to show that  $\sum_{\varphi^t} Q(\varphi^t, \alpha^t, \beta^t) = p(\alpha^t, \beta^t)$  by induction on  $t$ .



**Base case:**  $t = 1$ . Let us define  $Q(\varphi_1|\epsilon) = p(\varphi_1(\epsilon))$  and  $Q(\beta_1|\alpha^1, \epsilon) = p(\beta_1|\alpha_1)$ .

Then:

$$\begin{aligned}
\sum_{\varphi^1} Q(\varphi^1, \alpha^1, \beta^1) &= \sum_{\varphi_1} Q(\varphi_1, \alpha_1, \beta_1) \\
&= \sum_{\varphi_1} Q(\varphi_1|\epsilon, \epsilon, \epsilon) Q(\alpha_1|\varphi_1, \epsilon, \epsilon) Q(\beta_1|\varphi_1, \alpha_1, \epsilon) \quad (\text{by the chain rule}) \\
&= \sum_{\varphi_1} Q(\varphi_1|\epsilon) \delta_{\{\varphi_1(\epsilon)\}}(\alpha_1) Q(\beta_1|\alpha^1, \epsilon) \quad (\text{by Definition 6}) \\
&= \sum_{\varphi_1} p(\varphi_1(\epsilon)) \delta_{\{\varphi_1(\epsilon)\}}(\alpha_1) p(\beta_1|\alpha_1) \quad (\text{by construction of } Q) \\
&= p(\alpha_1) p(\beta_1|\alpha_1) \quad (\text{by definition of } \delta) \\
&= p(\alpha_1, \beta_1) \\
&= p(\alpha^1, \beta^1)
\end{aligned}$$

**Inductive case:** Let us define  $Q(\beta_t|\alpha^t, \beta^{t-1}) = p(\beta_t|\alpha^t, \beta^{t-1})$ , and

$$Q(\varphi_t|\varphi^{t-1}) = \prod_{\beta^{t-1}} p(\varphi_t(\beta^{t-1})|\varphi^{t-1}(\beta^{t-2}), \beta^{t-1})$$

Note that, if we consider  $\mathcal{X} = \{\beta^{t-1} \mid \beta_i \in \mathcal{B}, 1 \leq i \leq t-1\}$ ,  $\mathcal{Y} = \mathcal{A}$ , and  $p(\beta^{t-1}, \alpha_t) = p(\alpha_t|\varphi^{t-1}(\beta^{t-2}), \beta^{t-1})$ , then  $\mathcal{X}$ ,  $\mathcal{Y}$  and  $p$  satisfy the hypothesis of Lemma 2.

Then:

$$\begin{aligned}
&\sum_{\varphi^t} Q(\varphi^t, \alpha^t, \beta^t) \\
&= \quad (\text{by the chain rule}) \\
&\sum_{\varphi^t} Q(\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) Q(\varphi_t|\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) Q(\alpha_t|\varphi^t, \alpha^{t-1}, \beta^{t-1}) Q(\beta_t|\varphi^t, \alpha^t, \beta^{t-1}) \\
&= \quad (\text{by Definition 6}) \\
&\sum_{\varphi^t} Q(\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) Q(\varphi_t|\varphi^{t-1}) \delta_{\{\varphi_t(\beta^{t-1})\}}(\alpha_t) Q(\beta_t|\alpha^t, \beta^{t-1}) \\
&= \quad (\text{by construction of } Q) \\
&\sum_{\varphi^t} Q(\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) \left( \prod_{\beta^{t-1}} p(\varphi_t(\beta^{t-1})|\varphi^{t-1}(\beta^{t-2}), \beta^{t-1}) \right) \delta_{\{\varphi_t(\beta^{t-1})\}}(\alpha_t) p(\beta_t|\alpha^t, \beta^{t-1}) \\
&= \quad (\text{by definition of } \delta) \\
&\sum_{\substack{\varphi^t \\ \varphi_t(\beta^{t-1})=\alpha_t}} Q(\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) \left( \prod_{\beta^{t-1}} p(\varphi_t(\beta^{t-1})|\varphi^{t-1}(\beta^{t-2}), \beta^{t-1}) \right) p(\beta_t|\alpha^t, \beta^{t-1}) \\
&=
\end{aligned}$$

$$\begin{aligned}
& \sum_{\varphi^{t-1}} Q(\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) p(\beta_t | \alpha^t, \beta^{t-1}) \sum_{\varphi_t(\beta^{t-1}) = \alpha_t} \prod_{\beta^{t-1}} p(\varphi_t(\beta^{t-1}) | \varphi^{t-1}(\beta^{t-2}), \beta^{t-1}) \\
&= \text{(by Lemma 2)} \\
& \sum_{\varphi^{t-1}} Q(\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) \cdot p(\beta_t | \alpha^t, \beta^{t-1}) \cdot p(\alpha_t | \alpha^{t-1}, \beta^{t-1}) \\
&= \\
& p(\beta_t | \alpha^t, \beta^{t-1}) \cdot p(\alpha_t | \alpha^{t-1}, \beta^{t-1}) \cdot \sum_{\varphi^{t-1}} Q(\varphi^{t-1}, \alpha^{t-1}, \beta^{t-1}) \\
&= \text{(by induction hypothesis)} \\
& p(\beta_t | \alpha^t, \beta^{t-1}) \cdot p(\alpha_t | \alpha^{t-1}, \beta^{t-1}) \cdot p(\alpha^{t-1}, \beta^{t-1}) \\
&= \text{(by the chain rule)} \\
& p(\alpha^t, \beta^t)
\end{aligned}$$

□

**Corollary 1.** *Let a  $\mathcal{J}$  be a fully probabilistic IIHS. Let  $\{p(\beta_t | \alpha^t, \beta^{t-1})\}_{t=1}^T$  be a sequence of stochastic kernels and  $\{p(\alpha_t | \alpha^{t-1}, \beta^{t-1})\}_{t=1}^T$  a sequence of input distributions defined by  $\mathcal{J}$  according to Definitions 4 and 5. Then the reactor  $R = \{p(\varphi_t | \varphi^{t-1})\}_{t=1}^T$  compatible with respect to the  $\mathcal{J}$  is given by:*

$$p(\varphi_1) = p(\alpha_1 | \alpha^0, \beta^0) = p(\alpha_1) \quad (7)$$

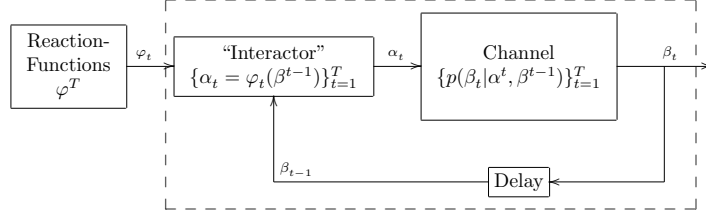
$$p(\varphi_t | \varphi^{t-1}) = \prod_{\beta^{t-1}} p(\varphi_t(\beta^{t-1}) | \varphi^{t-1}(\beta^{t-2}), \beta^{t-1}), \quad 2 \leq t \leq T \quad (8)$$

Figure 6 depicts the model for IIHS. Note that, in relation to Figure 2, there are some simplifications: (1) no message  $w$  is needed; (2) the encoder becomes an “interactor”; (3) the decoder is not used. At the beginning, a reaction function sequence  $\varphi^T$  is chosen and then the channel is used  $T$  times. At each usage  $t$ , the interactor produces the next input symbol  $\alpha_t$  by applying the reaction function  $\varphi_t$  to the fed back output  $\beta^{t-1}$ . Then the channel produces an output  $\beta_t$  based on the stochastic kernel  $p(\beta_t | \alpha^t, \beta^{t-1})$ . The output is then fed back to the encoder, which uses it for producing the next input.

We conclude this section by remarking on an intriguing coincidence: The notion of reaction function sequence  $\varphi^T$ , on the IIHSs, corresponds to the notion of deterministic scheduler [22]. In fact, each reaction function  $\varphi_t$  selects the next step,  $\alpha_t$ , on the basis of the  $\beta^{t-1}$  and  $\alpha^{t-1}$  (generated by  $\varphi^{t-1}$ ), and  $\beta^{t-1}, \alpha^{t-1}$  represent the path until that state.

## 5 Leakage in Interactive Systems

In this section we propose a definition for the notion of leakage in interactive systems. We first argue that mutual information is not the correct notion, and we propose to replace it with the directed information instead.



**Fig. 6.** Channel with memory and feedback model for IIHS

In the case of channels with memory and feedback, mutual information is defined as  $I(A^T; B^T) = H(A^T) - H(A^T|B^T)$ , and it is still symmetric (i.e.  $I(A^T; B^T) = I(B^T; A^T)$ ). Since the roles of  $A^T$  and  $B^T$  in  $I$  are interchangeable, this concept cannot capture *causality*, in the sense that it does not imply that  $A^T$  causes  $B^T$ , nor conversely. Mutual information expresses *correlation* between the sequences of random variables  $A^T$  and  $B^T$ .

Mathematically, for  $T$  usages of the channel, the mutual information  $I(A^T; B^T)$  can be expressed with the help of the chain rule of (2) in the following form.

$$I(A^T; B^T) = \sum_{t=1}^T I(A^t; B_t|B^{t-1}) \quad (9)$$

In the equation above, each term of the sum is the mutual information between the random variable  $B_t$  and the whole sequence of random variables  $A^T = A_1, \dots, A_T$ , given the history  $B^{t-1}$ . The equation emphasizes that at time  $1 \leq t \leq T$ , even though only the inputs  $\alpha^t = \alpha_1, \alpha_2, \dots, \alpha_t$  have been fed to the channel, the whole sequence  $A^T$ , including  $A_{t+1}, A_{t+2}, \dots, A_T$ , has a statistical correlation with  $B_t$ . Indeed, in the presence of feedback,  $B_t$  may influence  $A_{t+1}, A_{t+2}, \dots, A_T$ .

In order to show how the concept of directed information contrasts with the above, let us recall its definition:

$$I(A^T \rightarrow B^T) = \sum_{t=1}^T I(A^t; B_t|B^{t-1}).$$

$$I(B^T \rightarrow A^T) = \sum_{t=1}^T I(A_t; B^{t-1}|A^{t-1}).$$

These notions capture the concept of *causality*, to which the definition of mutual information is indifferent. The correlation between inputs and outputs  $I(A^T; B^T)$  is split into the information  $I(A^T \rightarrow B^T)$  that flows from input to output through the channel and the information  $I(B^T \rightarrow A^T)$  that flows from output to the input via feedback. Note that the directed information is not symmetric: the flow from  $A^T$  to  $B^T$  takes into account the correlation between  $A^t$  and  $B_t$ , while the flow from  $B^T$  to  $A^T$  takes into account the correlation between  $B^{t-1}$  and  $A_t$ .

It was proved in [29] that

$$I(A^T; B^T) = I(A^T \rightarrow B^T) + I(B^T \rightarrow A^T) \quad (10)$$

i.e., the mutual information is the sum of the directed information flow in both senses. Note that this formulation highlights the symmetry of mutual information from yet another perspective.

Once we split mutual information into directed information in the two opposite directions, it is important to understand the different role that the information flow in each direction plays.  $I(A^T \rightarrow B^T)$  represents the system behavior: via the channel the information flows from inputs to outputs according to the system specification, modeled by the channel stochastic kernels. This flow represents the amount of information an attacker can gain from the inputs by observing the outputs, and we argue that this is the real information leakage.

On the other hand,  $I(B^T \rightarrow A^T)$  represents how the environment reacts to the system: given the system outputs, the environment produces new inputs. We argue that the information flow from outputs to inputs is independent of any particular system: it is a characteristic of the environment itself. Hence, if an attacker knows how the environment reacts to outputs, i.e. the probabilistic behavior of the environment reactions given the system outputs, this knowledge is part of the *a priori* knowledge of the adversary. As a further justification, observe that this is a natural extension of the classical approach case, where the choice of secrets is seen as external to the system, i.e. determined by the environment. The probability distribution on the secrets constitutes the *a priori* knowledge and does not count as leakage. In order to encompass the classical approach, in our extended model we should preserve this principle, and a natural way to do so is to consider the secret choices, at every stage of the computation, as external. Their probability distributions, which are now in general conditional probability distributions depending on the history of secrets and observables, should therefore be considered as part of the external knowledge, and not counted as leakage.

The following example supports our claim that, in the presence of feedback, mutual information is not a correct notion of leakage.

*Example 5.* Consider the discrete memoryless channel with secret alphabet  $\mathcal{A} = \{a_1, a_2\}$  and observable alphabet  $\mathcal{B} = \{b_1, b_2\}$  whose matrix is represented in Table 4.

Suppose that the channel is used with feedback, in such a way that, for all  $1 \leq t \leq T$ , we have  $\alpha_{t+1} = a_1$  if  $\beta_t = b_1$ , and  $\alpha_{t+1} = a_2$  if  $\beta_t = b_2$ . It is easy to show that if  $T \geq 2$  then  $I(A^T; B^T) \neq 0$ . However, there is no leakage from  $A^T$  to  $B^T$ , since the rows of the matrix are all equal. We have indeed that  $I(A^T \rightarrow B^T) = 0$ , and the mutual information  $I(A^T; B^T)$  is only due to the feedback information flow  $I(B^T \rightarrow A^T)$ .

	$b_1$	$b_2$
$a_1$	0.5	0.5
$a_2$	0.5	0.5

**Table 4.** Channel matrix for Example 5

Having in mind the above discussion, we now propose a notion of information flow based on our model. We follow the idea of defining leakage and maximum leakage using the concepts of mutual information and capacity, making the necessary adaptations.

As discussed in the introduction, in the non interactive case the definition of leakage as mutual information, for a single use of the channel, is

$$I(A; B) = H(A) - H(A|B)$$

(cfr. for instance [4, 15]). This corresponds to view the leakage as difference between the a priori invulnerability degree,  $H(A)$ , and the a posteriori one,  $H(A|B)$ . The model of attacker which induces an invulnerability degree corresponding to Shannon entropy is discussed by Köpf and Basin in [15].

In the interactive case, we can extend this notion by considering the leakage at every step  $t$  as given by

$$I(A^t; B_t|B^{t-1}) = H(A^t|B^{t-1}) - H(A^t|B_t, B^{t-1})$$

The notion of attack is the same modulo the fact that we consider all the input from the beginning till step  $t$ , and the difference in its vulnerability induced by the observation of  $B_t$  (the output at step  $t$ ), taking into account the observation history  $B^{t-1}$ . It is then natural to consider as total leakage the summation of the contributions  $I(A^t; B_t|B^{t-1})$  for all the steps  $t$ . This is exactly the notion of directed information (cfr. Definition 2):

$$I(B^T \rightarrow A^T) = \sum_{t=1}^T I(A^t; B_t|B^{t-1})$$

**Definition 8.** *The information leakage of a fully probabilistic IHS is defined as the directed information  $I(A^T \rightarrow B^T)$  of the associated channel with memory and feedback.*

We now show an equivalent formulation of directed information that brings to a new interpretation in terms of attack model. First we need the following lemma.

**Lemma 3.**  $I(B^T \rightarrow A^T) = H(A^T) - \sum_{t=1}^T H(A_t|A^{t-1}, B^{t-1})$

*Proof.*

$$\begin{aligned} I(B^T \rightarrow A^T) &= \sum_{t=1}^T I(A_t; B^{t-1}|A^{t-1}) && \text{(by Definition 2)} \\ &= \sum_{t=1}^T (H(A_t|A^{t-1}) \\ &\quad - H(A_t|A^{t-1}, B^{t-1})) && \text{(by definition of mutual info.)} \\ &= H(A^T) - \sum_{t=1}^T H(A_t|A^{t-1}, B^{t-1}) && \text{(by the chain rule)} \end{aligned}$$

□

Next proposition points out the announced alternative formulation of directed information from input to output:

**Proposition 2.**  $I(A^T \rightarrow B^T) = \sum_{t=1}^T H(A_t|A^{t-1}, B^{t-1}) - H(A^T|B^T)$

*Proof.*

$$\begin{aligned}
I(A^T \rightarrow B^T) &= I(A^T; B^T) - I(B^T \rightarrow A^T) && \text{(by (10))} \\
&= I(A^T; B^T) - H(A^T) \\
&\quad + \sum_{t=1}^T H(A_t|A^{t-1}, B^{t-1}) && \text{(by Lemma 3)} \\
&= H(A^T) - H(A^T|B^T) - H(A^T) \\
&\quad + \sum_{t=1}^T H(A_t|A^{t-1}, B^{t-1}) && \text{(by definition of mutual info.)} \\
&= \sum_{t=1}^T H(A_t|A^{t-1}, B^{t-1}) - H(A^T|B^T)
\end{aligned}$$

□

We note that the term  $\sum_{t=1}^T H(A_t|A^{t-1}, B^{t-1})$  can be seen as the entropy  $H_R$  of the reactor  $R$ , i.e. the entropy of the inputs, taking into account their dependency on the previous outputs. This brings to an intriguing alternative interpretation of leakage:

*Remark 1.* The leakage can be seen as the difference between the a priori invulnerability degree of the whole secret  $A^T$ , assuming that the attacker knows the distribution of the reactor, and the a posteriori invulnerability degree, after the adversary has observed the whole output  $B^T$ .

In Section 6 we give an extensive and detailed example of how to calculate the leakage for a real security protocol.

In the case of secret-nondeterministic IHS, we have a stochastic kernel but no distribution on the code functions. In this case it seems natural to consider the worst leakage over all possible distributions on code functions. This is exactly the concept of capacity.

**Definition 9.** *The maximum leakage of a secret-nondeterministic IHS is defined as the capacity  $C_T$  of the associated channel with memory and feedback (cfr. (5)).*

A comparison with the definition of Gray (cfr. [13], Definition 5.3) is in order. As explained in the introduction, Gray's model is more complicated than ours, because it assumes that low and high variables are present at both ends of the channel. If we restrict the definition of Gray's capacity  $C^G$  to our case, by eliminating the low input and the high output, we obtain the following formula:

$$C_T^G = \sup_{\mathcal{D}_T} \frac{1}{T} \sum_{t=1}^T I(A^{t-1}; B_t|B^{t-1}) \quad (11)$$

By examining (11) against (9) and Definition 2, we can see that the only difference is that (11) considers the correlation between  $B_t$  and  $A^{t-1}$  instead than  $A^t$ . This seems to be intentional (cfr. [13], discussion after Definition 4.1). We are not sure why  $C^G$  is defined in this way, our best guess is that the high values must be those of the previous time step in order to encompass the theory of McLean [18]. In any case, Gray's conjecture that  $C_T^G$  corresponds to the channel transmission rate does not hold. For instance, it is easy to see that for  $T = 1$  we always have  $C_T^G = 0$ , but there obviously are channels which can transmit a non-zero amount of information even with one single use.

We conclude this section by showing that our approach to the notion of leakage generalizes the classical approach (based on mutual information) to the case of feedback. The idea is that, if a channel does not have feedback, then  $I(B^T \rightarrow A^T) = 0$  and therefore  $I(A^T; B^T) = I(A^T \rightarrow B^T)$ . In our opinion, the fact that mutual information turns out to be a particular case of directed information helps justifying the former as a good measure of information flow, despite its symmetry: in channels without feedback it is a good measure *because it coincides with directed information* from input to output.

**Lemma 4.** *In absence of feedback,  $I(B^T \rightarrow A^T) = 0$*

*Proof.* When feedback is not allowed,  $B^t$  and  $A_t$  are independent for  $1 \leq t \leq T$ . Then:

$$\begin{aligned}
 I(B^T \rightarrow A^T) &= \sum_{t=1}^T I(A_t; B^{t-1} | A^{t-1}) && \text{(by Definition 2)} \\
 &= \sum_{t=1}^T (H(A_t | A^{t-1}) - H(A_t | A^{t-1}, B^{t-1})) && \text{(by definition of mutual information)} \\
 &= \sum_{t=1}^T (H(A_t | A^{t-1}) - H(A_t | A^{t-1})) && \text{(by the independence of } B^{t-1} \text{ and } A^t) \\
 &= 0
 \end{aligned}$$

□

**Proposition 3.** *In absence of feedback, leakage can be equivalently defined as directed information or as mutual information. Similarly, in absence of feedback, the maximum leakage can be equivalently defined as directed capacity or as capacity.*

*Proof.* It follows directly from Lemma 4 and (10). □

## 6 Modeling IIHSs as channels: An example

In this section we show the application of our approach to the *Cocaine Auction Protocol* [25]. The formalization of this protocol in terms of IIHSs using our framework makes it possible to prove the claim in [25] suggesting that if the seller knows the identity of the bidders then the (strong) anonymity guaranties are not provided anymore.

Let us consider a scenario in which several mobsters are gathered around a table. An auction is about to be held in which one of them offers his next shipment of cocaine to the highest bidder. The seller describes the merchandise and proposes a starting

price. The others then bid increasing amounts until there are no bids for 30 consecutive seconds. At that point the seller declares the auction closed and arranges a secret appointment with the winner to deliver the goods.

The basic protocol is fairly simple and is organized as a succession of rounds of bidding. Round  $i$  starts with the seller announcing the bid price  $b_i$  for that round. Buyers have  $t$  seconds to make an offer (i.e. to say yes, meaning “I’m willing to buy at the current bid price  $b_i$ ”). As soon as one buyer anonymously says yes, he becomes the winner  $w_i$  of that round and a new round begins. If nobody says anything for  $t$  seconds, round  $i$  is concluded by timeout and the auction is won by the winner  $w_{i-1}$  of the previous round, if one exists. If the timeout occurs during round 0, this means that nobody made any offers at the initial price  $b_0$ , so there is no sale.

Although our framework allows the formalization of this protocol for an arbitrary number of bidders and bidding rounds, for illustration purposes we will consider the case of two bidders (*Candlemaker* and *Scarface*) and two rounds of bids. Furthermore, we assume that the initial bid is always 1 dollar, so the first bid does not need to be announced by the seller. In each turn the seller can choose how much he wants to increase the current bid value. This is done by adding an increment to the last bid. There are two options of increments, namely  $inc_1$  (1 dollar) and  $inc_2$  (2 dollars). In that way,  $b_{i+1}$  is either  $b_i + inc_1$  or  $b_i + inc_2$ . We can describe this protocol as a *normalized* IHHS  $\mathcal{I} = (M, \mathcal{A}, \mathcal{B}, \mathcal{C})$ , where  $\mathcal{A} = \{\text{Candlemaker}, \text{Scarface}, a_*\}$  is the set of secret actions,  $\mathcal{B} = \{inc_1, inc_2, b_*\}$  is the set of observable actions,  $\mathcal{C} = \emptyset$  is the set of hidden actions, and the probabilistic automaton  $M$  is represented in Figure 7. For clarity reasons, transitions with probability 0 are not represented in the automaton. Note that the special secret action  $a_*$  represents the situation where neither *Candlemaker* nor *Scarface* bid. The special observable action  $b_*$  is only possible after no one has bidden, and signals the end of the auction and, therefore, no further bids are allowed.

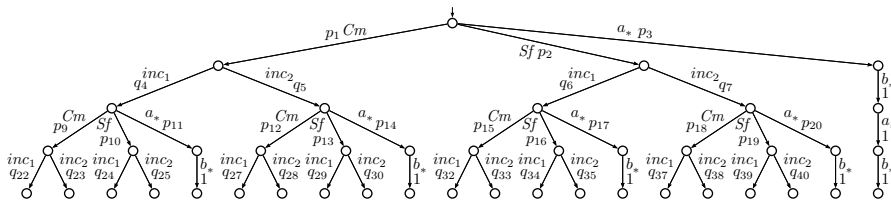


Fig. 7. Cocaine Auction example

Table 5 shows all the stochastic kernels for this example.

The interested reader can find the construction of the reaction functions in the Appendix.



(a) $t=1, p(\beta_1 \alpha^1, \beta^0)$				(b) $t=2, p(\beta_2 \alpha^2, \beta^1)$			
$\alpha_1 \rightarrow \beta_1$	$inc_1$	$inc_2$	$b_*$	$\alpha_1, \beta_1, \alpha_2 \rightarrow \beta_2$	$inc_1$	$inc_2$	$b_*$
<i>Candlemaker</i>	$q_4$	$q_5$	0	<i>Candlemaker, inc<sub>1</sub>, Candlemaker</i>	$q_{22}$	$q_{23}$	0
<i>Scarface</i>	$q_6$	$q_7$	0	<i>Candlemaker, inc<sub>1</sub>, Scarface</i>	$q_{24}$	$q_{25}$	0
$a^*$	0	0	1	<i>Candlemaker, inc<sub>1</sub>, a<sub>*</sub></i>	0	0	1
				<i>Candlemaker, inc<sub>2</sub>, Candlemaker</i>	$q_{27}$	$q_{28}$	0
				<i>Candlemaker, inc<sub>2</sub>, Scarface</i>	$q_{29}$	$q_{30}$	0
				<i>Candlemaker, inc<sub>2</sub>, a<sub>*</sub></i>	0	0	1
				<i>Scarface, inc<sub>1</sub>, Candlemaker</i>	$q_{32}$	$q_{33}$	0
				<i>Scarface, inc<sub>1</sub>, Scarface</i>	$q_{34}$	$q_{35}$	0
				<i>Scarface, inc<sub>1</sub>, a<sub>*</sub></i>	0	0	1
				<i>Scarface, inc<sub>2</sub>, Candlemaker</i>	$q_{37}$	$q_{38}$	0
				<i>Scarface, inc<sub>2</sub>, Scarface</i>	$q_{39}$	$q_{40}$	0
				<i>Scarface, inc<sub>2</sub>, a<sub>*</sub></i>	0	0	1
				$a_*, b_*, a_*$	0	0	1
				All other lines	0	0	1

**Table 5.** Stochastic kernels for the Cocaine Auction example.

### 6.1 Calculating the information leakage

Let us now calculate the information leakage for this example using the concepts from Section 5. We are going to analyze three different scenarios:

**Example a:** There is feedback, but the probability of an observable does not depend on the history of secrets. In the auction protocol, this corresponds to a scenario where the probability of one of the mobsters to bid can depend on the increment imposed by the seller, but the history of who has previously bid in the past has no influence on how the seller chooses the bid increment in the coming turns. In other words, the seller cannot use the information of who has been bidding to change his strategy of defining the new increments. This situation corresponds to the original description of the protocol in [25], where the seller does not have access to the identity of the bidder, for the sake of anonymity preservation. In general, we have  $p(\beta_t|\alpha^t, \beta^{t-1}) = p(\beta_t|\beta^{t-1})$  for every  $1 \leq t \leq T$ . However, there is an exception: if there is no bidder, the case modeled by the secret being  $a_*$ , then the auction terminates, which is signaled by the observable  $b_*$ .

**Example b:** This is the most general case, without any restrictions. The presence of feedback allows the probability of the bidder to depend of the increment on the price. For instance, if *Candlemaker* is richer than *Scarface*, it is more likely that the former bids if the increment in the price is  $inc_2$  instead of  $inc_1$ . Also, the probability of an observable can depend on the history of secrets, i.e., in general  $p(\beta_t|\alpha^t, \beta^{t-1}) \neq p(\beta_t|\beta^{t-1})$  for  $1 \leq t \leq T$ . This scenario can represent a situation where the seller is corrupted and can use his information to affect the outcome of the auction. As an example, suppose that the seller is a friend of *Scarface* and he wants to help him in the auction. One way of doing so is to check who was the winner of the last bidding round. Whenever the winner is *Candlemaker*, the seller chooses as increment the small value  $inc_1$ , hoping that it will give *Scarface* a good chance to bid in the next round. On the other hand, whenever the seller detects that

the winner is *Scarface*, he chooses as the next increment the greater value  $inc_2$ , hoping that it will minimize the chances of *Candlemaker* to bid in the next round (and therefore maximizing the chances of the auction to end up having *Scarface* as the final winner).

**Example c:** There is no feedback. In the cocaine auction, we can have the (maybe unrealistic) situation in which the increment added to the bid has no influence on the probability of *Candlemaker* or *Scarface* being the bidder. Mathematically, we have  $p(\alpha_t|\alpha^{t-1}, \beta^{t-1}) = p(\alpha_t|\alpha^{t-1})$  for every  $1 \leq t \leq T$ . However, as in Example b, we do not impose any restriction to  $p(\beta_t|\alpha^t, \beta^{t-1})$ .

For each scenario we need to fill in the values of the probabilities in the protocol tree in Figure 7. The probabilities for each example are listed in Table 6.

Table 7 shows a comparison between some relevant values on the three cases.

In Example a, since the probability of observables does not depend on the history of secrets, there is (almost) no information flowing from the input to the output, and the directed information  $I(A^T \rightarrow B^T)$  is close to zero, i.e., the leakage is low. The only reason why the leakage is not zero is because the end of an auction needs to be signaled. However, due to presence of feedback, the directed information in the other sense  $I(B^T \rightarrow A^T)$  is non-zero, and so is the mutual information  $I(A^T; B^T)$ . This is an example where the mutual information does not correspond to the real information leakage, since some (in this case, most) of the correlation between input and output can be attributed to the feedback.

In Example b the information flow from input to output  $I(A^T \rightarrow B^T)$  is significantly higher than zero, but still, due to feedback, the information flow from outputs to inputs  $I(B^T \rightarrow A^T)$  is not zero and the mutual information  $I(A^T; B^T)$  is higher than the directed information  $I(A^T \rightarrow B^T)$ .

In Example c, the absence of feedback implies that  $I(B^T \rightarrow A^T)$  is zero. In that case the values of  $I(A^T; B^T)$  and  $I(A^T \rightarrow B^T)$  coincide, and represent the real leakage.

## 7 Topological properties of IIHSs and their Capacity

In this section we show how to extend to IIHSs the notion of pseudometric defined in [9] for Concurrent Labelled Markov Chains, and we prove that the capacity of the corresponding channels is a continuous function with respect to this pseudometric. The pseudometric construction is sound for general IIHSs, but the result on capacity is only valid for secret-nondeterministic IIHSs.

Given a set of states  $S$ , a pseudometric is a function  $d$  that yields a non-negative real number for each pair of states and satisfies the following:  $d(s, s) = 0$ ;  $d(s, t) = d(t, s)$ , and  $d(s, t) \leq d(s, u) + d(u, t)$ . We say that a pseudometric  $d$  is  $c$ -bounded if  $\forall s, t : d(s, t) \leq c$ , where  $c$  is a positive real number.

Note that, in contrast to metrics, in pseudometrics two elements can have distance 0 without being identical. The reason for considering pseudometrics instead than metrics is because the purpose is to extend the notion of (probabilistic) bisimulation: having distance 0 will correspond to being bisimilar.

Probability variable	Example a value	Example b value	Example c value
$p_1$	0.7	0.7	0.7
$p_2$	0.2	0.2	0.2
$p_3$	0.1	0.1	0.1
$q_4$	0.9	0.1	0.1
$q_5$	0.1	0.9	0.9
$q_6$	0.9	0.9	0.9
$q_7$	0.1	0.1	0.1
$p_9$	0.6	0.6	0.6
$p_{10}$	0.3	0.3	0.3
$p_{11}$	0.1	0.1	0.1
$p_{12}$	0.5	0.5	0.6
$p_{13}$	0.3	0.3	0.3
$p_{14}$	0.2	0.2	0.1
$p_{15}$	0.4	0.4	0.5
$p_{16}$	0.4	0.4	0.2
$p_{17}$	0.2	0.2	0.3
$p_{18}$	0.6	0.6	0.5
$p_{19}$	0.3	0.3	0.2
$p_{20}$	0.1	0.1	0.3
$q_{22}$	0.4	0.1	0.1
$q_{23}$	0.6	0.9	0.9
$q_{24}$	0.7	0.9	0.9
$q_{25}$	0.3	0.1	0.1
$q_{27}$	0.2	0.1	0.1
$q_{28}$	0.8	0.9	0.9
$q_{29}$	0.1	0.9	0.9
$q_{30}$	0.9	0.1	0.1
$q_{32}$	0.4	0.1	0.1
$q_{33}$	0.6	0.9	0.9
$q_{34}$	0.7	0.9	0.9
$q_{35}$	0.3	0.1	0.1
$q_{37}$	0.2	0.1	0.1
$q_{38}$	0.8	0.9	0.9
$q_{39}$	0.1	0.9	0.9
$q_{40}$	0.9	0.1	0.1

**Table 6.** Values of the probabilities in Figure 7 in 3 different examples.

We now define a complete lattice on pseudometrics, in order to define the distance between IIHSs as the greatest fixpoint of a particular transformation, in line with the coinductive theory of bisimilarity. Since larger bisimulations identify more, the natural extension of the ordering of pseudometrics must shorten the distances as we go up in the lattice:

Interpretation	Symbol	Example a	Example b	Example c
Input uncertainty	$H(A^T)$	2.3833	2.4891	2.3607
Reactor uncertainty	$H_R$	2.3768	2.4832	2.3607
A posteriori uncertainty	$H(A^T B^T)$	1.3683	0.0677	0.6646
Mutual information	$I(A^T; B^T) = H(A^T) - H(A^T B^T)$	1.0150	1.8214	1.6961
Leakage	$I(A^T \rightarrow B^T) = H_R - H(A^T B^T)$	1.0085	1.8155	1.6961
Feedback information	$I(B^T \rightarrow A^T)$	0.185955	0.0060	0.0000

**Table 7.** Values for the examples.

**Definition 10.**  $\mathcal{M}$  is the class of 1-bounded pseudometrics on states with the ordering

$$d \preceq d' \text{ if } \forall s, s' \in S : d(s, s') \geq d'(s, s').$$

It is easy to see that  $(\mathcal{M}, \preceq)$  is a complete lattice. In order to define pseudometrics on IIHSs, we now need to lift the pseudometrics on states to pseudometrics on distributions in  $\mathcal{D}(\mathcal{L} \times S)$ . Following standard lines [30, 9, 8], we apply the construction based on the Kantorovich metric [14].

**Definition 11.** For  $d \in \mathcal{M}$ , and  $\mu, \mu' \in \mathcal{D}(\mathcal{L} \times S)$ , we define  $d(\mu, \mu')$  (overloading the notation  $d$ ) as  $d(\mu, \mu') = \max \sum_{(\ell_i, s_i) \in \mathcal{L} \times S} (\mu(\ell_i, s_i) - \mu'(\ell_i, s_i)) x_i$  where the maxima is taken over all possible values of the  $x_i$ 's, subject to the constraints  $0 \leq x_i \leq 1$  and  $x_i - x_j \leq \hat{d}((\ell_i, s_i), (\ell_j, s_j))$ , where

$$\hat{d}((\ell_i, s_i), (\ell_j, s_j)) = \begin{cases} 1 & \text{if } \ell_i \neq \ell_j \\ d(s_i, s_j) & \text{otherwise} \end{cases}$$

It can be shown that with this definition  $m$  is a pseudometric on  $\mathcal{D}(\mathcal{L} \times S)$ .

**Definition 12.** A pseudometric  $d \in \mathcal{M}$  is a bisimulation pseudometric<sup>3</sup> if, for all  $\epsilon \in [0, 1)$ ,  $d(s, s') \leq \epsilon$  implies that if  $s \rightarrow \mu$ , then there exists some  $\mu'$  such that  $s' \rightarrow \mu'$  and  $d(\mu, \mu') \leq \epsilon$ .

Note that it is not necessary to require the converse of the condition in Definition 12 to get a complete analogy with bisimulation: the converse is indeed implied by the symmetry of  $d$  as a pseudometric. Note also that we prohibit  $\epsilon$  to be 1 because throughout this paper 1 represents the maximum distance, which includes the case where one state may perform a transition and the other may not.

The greatest bisimulation pseudometric is

$$d_{max} = \bigsqcup \{d \in \mathcal{M} \mid d \text{ is a bisimulation pseudometric}\}$$

We now characterize  $d_{max}$  as a fixed point of a monotonic function  $\Phi$  on  $\mathcal{M}$ . Eventually we are interested in the distance between IIHSs, and for the sake of simplicity, from

<sup>3</sup> In literature a pseudometric with this property is also known as bisimulation metric, although it is still a pseudometric.

now on we consider only the distance between states belonging to different IIHSs. The extension to the general case is trivial. For clarity purposes, we assume that different IIHSs have disjoint sets of states.

**Definition 13.** Given two IIHSs with transition relations  $\theta$  and  $\theta'$  respectively, and a pseudometric  $d$  on states, define  $\Phi : \mathcal{M} \rightarrow \mathcal{M}$  as:

$$\Phi(d)(s, s') = \begin{cases} \max_i d(s_i, s'_i) & \text{if } \vartheta(s) = \{\delta_{(a_1, s_1)}, \dots, \delta_{(a_m, s_m)}\} \\ & \text{and } \vartheta'(s') = \{\delta_{(a_1, s'_1)}, \dots, \delta_{(a_m, s'_m)}\} \\ d(\mu, \mu') & \text{if } \vartheta(s) = \{\mu\} \text{ and } \vartheta'(s') = \{\mu'\} \\ 0 & \text{if } \vartheta(s) = \vartheta'(s') = \emptyset \\ 1 & \text{otherwise} \end{cases}$$

It is easy to see that the definition of  $\Phi$  is a particular case of the function  $F$  defined in [9, 8]. Hence it can be proved, by adapting the proofs of the analogous results in [9, 8], that  $F(d)$  is a pseudometric, and that the following property holds.

**Lemma 5.** For  $\epsilon \in [0, 1)$ ,  $\Phi(d)(s, s') \leq \epsilon$  holds if and only if whenever  $s \rightarrow \mu$ , there exists some  $\mu'$  such that  $s' \rightarrow \mu'$  and  $d(\mu, \mu') \leq \epsilon$ .

**Corollary 2.** A pseudometric  $d$  is a bisimulation pseudometric iff  $d \preceq \Phi(d)$ .

As a consequence of Corollary 2, we have that  $d_{max} = \bigsqcup \{d \in \mathcal{M} \mid d \preceq \Phi(d)\}$ , and still as a particular case of  $F$  in [9, 8], we have that  $\Phi$  is monotonic on  $\mathcal{M}$ .

We can now apply Tarski's fixed point theorem, which ensures that  $d_{max}$  is the greatest fixed point of  $\Phi$ . Furthermore, by Corollary 2 we know that  $d_{max}$  is indeed a bisimulation pseudometric, and that it is the greatest bisimulation pseudometric. In addition, the finite branching property of IIHSs ensures that the closure ordinal of  $\Phi$  is  $\omega$  (cf. Lemma 3.10 in the full version of [9], available on the authors' web pages). Therefore one can proceed in a standard way to show that  $d_{max} = \prod \{\Phi^i(\top) \mid i \in \mathbb{N}\}$ , where  $\top$  is the greatest pseudometric (i.e.  $\top(s, s') = 0$  for every  $s, s'$ ), and  $\Phi^0(\top) = \top$ .

Given two IIHSs  $\mathcal{J}$  and  $\mathcal{J}'$ , with initial states  $s$  and  $s'$  respectively, we define the distance between  $\mathcal{J}$  and  $\mathcal{J}'$  as  $d(\mathcal{J}, \mathcal{J}') = d_{max}(s, s')$ . The following properties are auxiliary to the theorem which states the continuity of the capacity.

**Lemma 6.** Consider two IIHSs  $\mathcal{J}$  and  $\mathcal{J}'$  with transition functions  $\vartheta$  and  $\vartheta'$  respectively. Given  $t \geq 2$  and two sequences  $\alpha^t$  and  $\beta^t$ , assume that both  $\mathcal{J}(\alpha^{t-1}, \beta^{t-1})$  and  $\mathcal{J}'(\alpha^{t-1}, \beta^{t-1})$  are defined, that  $d_{max}(\mathcal{J}(\alpha^{t-1}, \beta^{t-1}), \mathcal{J}'(\alpha^{t-1}, \beta^{t-1})) < p(\beta_t \mid \alpha^t, \beta^{t-1})$ , and  $\vartheta(\mathcal{J}(\alpha^t, \beta^{t-1})) \neq \emptyset$ . Then:

1.  $\vartheta'(\mathcal{J}'(\alpha^t, \beta^{t-1})) \neq \emptyset$  holds as well,
2.  $\mathcal{J}(\alpha^t, \beta^t)$  and  $\mathcal{J}'(\alpha^t, \beta^t)$  are both defined,  $p(\beta_t \mid \alpha^t, \beta^{t-1}) > 0$ , and

$$d_{max}(\mathcal{J}(\alpha^t, \beta^t), \mathcal{J}'(\alpha^t, \beta^t)) \leq \frac{d_{max}(\mathcal{J}(\alpha^{t-1}, \beta^{t-1}), \mathcal{J}'(\alpha^{t-1}, \beta^{t-1}))}{p(\beta_t \mid \alpha^t, \beta^{t-1})}$$

*Proof.*

1. Assume  $\vartheta(\mathcal{J}(\alpha^t, \beta^{t-1})) \neq \emptyset$  and, by contradiction,  $\vartheta'(\mathcal{J}'(\alpha^t, \beta^{t-1})) = \emptyset$ . Since  $d_{max}$  is a fixed point of  $F$ , we have  $d_{max} = F(d_{max})$ , and therefore

$$\begin{aligned} d_{max}(\mathcal{J}(\alpha^t, \beta^{t-1}), \mathcal{J}'(\alpha^t, \beta^{t-1})) &= F(d_{max})(\mathcal{J}(\alpha^t, \beta^{t-1}), \mathcal{J}'(\alpha^t, \beta^{t-1})) \\ &= 1 \\ &\geq p(\beta_t | \alpha^t, \beta^{t-1}), \end{aligned}$$

against the hypothesis.

2. If  $\vartheta(\mathcal{J}(\alpha^t, \beta^{t-1})) \neq \emptyset$ , then, by the first point of this lemma,  $\vartheta'(\mathcal{J}'(\alpha^t, \beta^{t-1})) \neq \emptyset$  holds as well, and therefore both  $\mathcal{J}(\alpha^t, \beta^t)$  and  $\mathcal{J}'(\alpha^t, \beta^t)$  are defined. The hypothesis  $d_{max}(\mathcal{J}(\alpha^{t-1}, \beta^{t-1}), \mathcal{J}'(\alpha^{t-1}, \beta^{t-1})) < p(\beta_t | \alpha^t, \beta^{t-1})$  ensures that  $p(\beta_t | \alpha^t, \beta^{t-1}) < 0$ . Let us now prove the bound on  $d_{max}(\mathcal{J}(\alpha^t, \beta^t), \mathcal{J}'(\alpha^t, \beta^t))$ . By definition of  $\Phi$ , we have

$$\Phi(d_{max})(\mathcal{J}(\alpha^{t-1}, \beta^{t-1}), \mathcal{J}'(\alpha^{t-1}, \beta^{t-1})) \geq d_{max}(\mathcal{J}(\alpha^t, \beta^{t-1}), \mathcal{J}'(\alpha^t, \beta^{t-1})).$$

Since  $d_{max} = \Phi(d_{max})$ , we have

$$d_{max}(\mathcal{J}(\alpha^{t-1}, \beta^{t-1}), \mathcal{J}'(\alpha^{t-1}, \beta^{t-1})) \geq d_{max}(\mathcal{J}(\alpha^t, \beta^{t-1}), \mathcal{J}'(\alpha^t, \beta^{t-1})). \quad (12)$$

By definition of  $\Phi$  and of the Kantorovich metric, we have

$$\begin{aligned} \Phi(d_{max})(\mathcal{J}(\alpha^t, \beta^{t-1}), \mathcal{J}'(\alpha^t, \beta^{t-1})) &\geq p(\beta_t | \alpha^t, \beta^{t-1}) \cdot \\ &\quad d_{max}(\mathcal{J}(\alpha^t, \beta^t), \mathcal{J}'(\alpha^t, \beta^t)). \end{aligned}$$

Using again  $d_{max} = \Phi(d_{max})$ , we get

$$\begin{aligned} d_{max}(\mathcal{J}(\alpha^t, \beta^{t-1}), \mathcal{J}'(\alpha^t, \beta^{t-1})) &\geq p(\beta_t | \alpha^t, \beta^{t-1}) \cdot \\ &\quad d_{max}(\mathcal{J}(\alpha^t, \beta^t), \mathcal{J}'(\alpha^t, \beta^t)), \end{aligned}$$

which, together with (12), allows us to conclude. □

**Lemma 7.** Consider two IHSSs  $\mathcal{J}$  and  $\mathcal{J}'$ , and let  $p(\cdot | \cdot, \cdot)$  and  $p'(\cdot | \cdot, \cdot)$  be their distributions on the output nodes. Given  $T > 0$ , and two sequences  $\alpha^T$  and  $\beta^T$ , assume that  $p(\beta_t | \alpha^t, \beta^{t-1}) > 0$  for every  $t < T$ . Let  $m = \min_{1 \leq t < T} p(\beta_t | \alpha^t, \beta^{t-1})$  and let  $\epsilon \in (0, m^{T-1})$ . Assume  $d(\mathcal{J}, \mathcal{J}') < \epsilon$ . Then, for every  $t \leq T$ , we have

$$p(\beta_t | \alpha^t, \beta^{t-1}) - p'(\beta_t | \alpha^t, \beta^{t-1}) < \frac{\epsilon}{m^{T-1}}.$$

*Proof.* Observe that, for every  $t < T$ ,  $\mathcal{J}(\alpha^t, \beta^t)$  must be defined, and, by repeatedly applying Lemma 6(1), we get that also  $\mathcal{J}'(\alpha^t, \beta^t)$  is defined. By definition of  $\varphi$ , and of the Kantorovich metric, we have

$$p(\beta_t | \alpha^t, \beta^{t-1}) - p'(\beta_t | \alpha^t, \beta^{t-1}) \leq \Phi(d_{max})(\mathcal{J}(\alpha^{t-1}, \beta^{t-1}), \mathcal{J}'(\alpha^{t-1}, \beta^{t-1})),$$

and since  $d_{max}$  is a fixed point of  $\Phi$ , we get

$$p(\beta_t | \alpha^t, \beta^{t-1}) - p'(\beta_t | \alpha^t, \beta^{t-1}) \leq d_{max}(\mathcal{J}(\alpha^{t-1}, \beta^{t-1}), \mathcal{J}'(\alpha^{t-1}, \beta^{t-1})). \quad (13)$$

By applying  $t - 1$  times Lemma 6(2), from (13) we get

$$\begin{aligned} p(\beta_t | \alpha^t, \beta^{t-1}) - p'(\beta_t | \alpha^t, \beta^{t-1}) &\leq \frac{d_{max}(\mathcal{J}(\alpha^0, \beta^0), \mathcal{J}'(\alpha^0, \beta^0))}{m^{t-1}} \\ &= \frac{d(\mathcal{J}, \mathcal{J}')}{m^{t-1}} \\ &\leq \frac{d(\mathcal{J}, \mathcal{J}')}{m^{T-1}} \\ &< \frac{\epsilon}{m^{T-1}} \end{aligned}$$

□

Note that previous lemma states a sort of continuity property of the matrices obtained from IIHSs, but not uniform continuity, because of the dependence on one of the two IIHSs. It is easy to see (from the proof of the Lemma) that uniform continuity does not hold.

The main contribution of this section, stated in next theorem, is the continuity of the capacity w.r.t. the pseudometric on IIHSs. For this theorem, we assume that the IIHSs are normalized. Furthermore, it is crucial that they are secret-nondeterministic (while the definition of the pseudometric holds in general).

**Theorem 4.** *Consider two normalized IIHSs  $\mathcal{J}$  and  $\mathcal{J}'$ , and fix a  $T > 0$ . For every  $\epsilon > 0$  there exists  $\nu > 0$  such that if  $d(\mathcal{J}, \mathcal{J}') < \nu$  then  $|C_T(\mathcal{J}) - C_T(\mathcal{J}')| < \epsilon$ .*

*Proof.* Consider two normalized IIHSs  $\mathcal{J}$  and  $\mathcal{J}'$  and choose  $T, \epsilon > 0$ . Observe that

$$\begin{aligned} |C_T(\mathcal{J}) - C_T(\mathcal{J}')| &= \left| \max_{p_F(\cdot)} \frac{1}{T} I(A^T \rightarrow B^T) - \max_{p_F(\cdot)} \frac{1}{T} I(A'^T \rightarrow B'^T) \right| \\ &\leq \frac{1}{T} \max_{p_F(\cdot)} |I(A^T \rightarrow B^T) - I(A'^T \rightarrow B'^T)| \end{aligned}$$

Since the directed information  $I(A^T \rightarrow B^T)$  is defined by means of arithmetic operations and logarithms on the joint probabilities  $p(\alpha^t, \beta^t)$  and on the conditional probabilities  $p(\alpha^t, \beta^t), p(\alpha^t, \beta^{t-1})$ , which in turn can be obtained by means of arithmetic operations from the probabilities  $p(\beta_t | \alpha^t, \beta^{t-1})$  and  $p_F(\varphi^t)$ , we have that  $I(A^T \rightarrow B^T)$  is a continuous functions of the distributions  $p(\beta_t | \alpha^t, \beta^{t-1})$  and  $p_F(\varphi^t)$ , for every  $t \leq T$ . Let  $p(\beta_t | \alpha^t, \beta^{t-1}), p'(\beta_t | \alpha^t, \beta^{t-1})$  be the distributions on the output nodes of  $\mathcal{J}$  and  $\mathcal{J}'$ , modified in the following way: starting from level  $T$ , whenever  $p(\beta_t | \alpha^t, \beta^{t-1}) = 0$ , then we redefine the distributions at all the output nodes of the subtree rooted in  $\mathcal{J}(\alpha^t, \beta^t)$  so that they coincide with the distribution of the corresponding nodes of in  $\mathcal{J}'$ , and analogously for  $p'(\beta_t | \alpha^t, \beta^{t-1})$ . Note that this transformation does not change the directed information, because the subtree rooted in  $\mathcal{J}(\alpha^t, \beta^t)$  does not contribute to it, due to the fact that it depends the probability of reaching any of its nodes is 0. The continuity of  $I(A^T \rightarrow B^T)$  implies that there exists  $\epsilon' > 0$  such that, if  $|p(\beta_t | \alpha^t, \beta^{t-1}) - p'(\beta_t | \alpha^t, \beta^{t-1})| < \epsilon'$  for all  $t \leq T$  and all sequences  $\alpha^t, \beta^t$ , then,

for any  $p_F(\varphi^t)$ , we have  $|I(A^T \rightarrow B^T) - I(A'^T \rightarrow B'^T)| < \epsilon$ . The result then follows from Lemma 7, by choosing

$$\nu = \epsilon' \cdot \min\left(\min_{\substack{1 \leq t < T \\ p(\beta_t | \alpha^t, \beta^{t-1}) > 0}} p(\beta_t | \alpha^t, \beta^{t-1}), \min_{\substack{1 \leq t < T \\ p'(\beta_t | \alpha^t, \beta^{t-1}) > 0}} p'(\beta_t | \alpha^t, \beta^{t-1})\right).$$

□

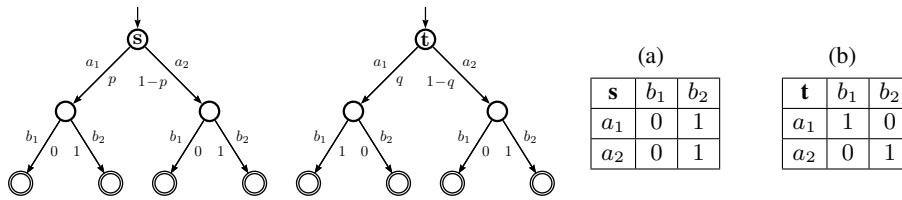
We conclude this section with an example showing that the continuity result for the capacity does not hold if the construction of the channel is done starting from a system in which the secrets are endowed with a probability distribution. This is also the reason why we could not simply adopt the proof technique of the continuity result in [9] and we had to come up with a different reasoning.

*Example 6.* Consider the two following programs, where  $a_1, a_2$  are secrets,  $b_1, b_2$  are observable,  $\parallel$  is the parallel operator, and  $+_p$  is a binary probabilistic choice that assigns probability  $p$  to the left branch, and probability  $1 - p$  to the right one.

- s)  $(\text{send}(a_1) +_p \text{send}(a_2)) \parallel \text{receive}(x).\text{output}(b_2)$
- t)  $(\text{send}(a_1) +_q \text{send}(a_2)) \parallel \text{receive}(x).\text{if } x = a_1 \text{ then output}(b_1) \text{ else output}(b_2)$ .

Table 8 shows the fully probabilistic IHSs corresponding to these programs, and their associated channels, which in this case (since the secret actions are all at the top-level) are classical channels, i.e. memoryless and without feedback. As usual for classic channels, they do not depend on  $p$  and  $q$ . It is easy to see that the capacity of the first channel is 0 and the capacity of the second one is 1. Hence their difference is 1, independently from  $p$  and  $q$ .

Let now  $p = 0$  and  $q = \epsilon$ . It is easy to see that the distance between  $s$  and  $t$  is  $\epsilon$ . Therefore (when the automata have probabilities on the secrets), the capacity is not a continuous function of the distance.



**Table 8.** The IHSs of Example 6 and their corresponding channels, (a) for  $s$  and (b) for  $t$ .



## 8 Conclusion and discussion

In this paper we have investigated the problem of information leakage in interactive systems, and proved that these systems can be modeled as channels with memory and feedback. The situation is summarized in Table 9(a). The comparison with the classical situation of non-interactive systems is represented in (b). Furthermore, we have proved that the channel capacity is a continuous function of a pseudometric based on the Kantorovich metric.

IIHSs as automata	IIHSs as channels	Notion of leakage
Normalized IIHSs with nondeterministic secrets and probabilistic observables	Sequence of stochastic kernels $\{p(\beta_t \alpha^t, \beta^{t-1})\}_{t=1}^T$	Leakage as capacity
Normalized IIHSs with a deterministic scheduler solving the nondeterminism	Sequence of stochastic kernels $\{p(\beta_t \alpha^t, \beta^{t-1})\}_{t=1}^T$ + reaction function seq. $\varphi^T$	
Fully probabilistic normalized IIHSs	Sequence of stochastic kernels $\{p(\beta_t \alpha^t, \beta^{t-1})\}_{t=1}^T$ + reactor $\{p(\varphi_t \varphi^{t-1})\}_{t=1}^T$	Leakage as directed information $I(A^T \rightarrow B^T)$

(a)

Classical channels	Channels with memory and feedback
The system is modeled in independent uses of the channel, often a unique use.	The system is modeled in several consecutive uses of the channel.
The channel is from $\mathcal{A}^T \rightarrow \mathcal{B}^T$ , i.e., its input is a single string $\alpha^T = \alpha_1 \dots \alpha_T$ of secret symbols and its output is a single string $\beta^T = \beta_1 \dots \beta_T$ of observable symbols.	The channel is from $\mathcal{F} \rightarrow \mathcal{B}$ , i.e. its input is a reaction function $\varphi_t$ and its output is an observable $\beta_t$ .
The channel is memoryless and in general implicitly it is assumed the absence of feedback.	The channel has memory. Despite the fact that the channel from $\mathcal{F} \rightarrow \mathcal{B}$ does not have feedback, the internal stochastic kernels do.
The capacity is calculated using information $I(A^T; B^T)$ .	The capacity is calculated using mutual directed information $I(A^T \rightarrow B^T)$ .

(b)

Table 9.

Throughout the paper we have assumed that the probability distributions on the secret choices are part of the external knowledge and, therefore, not considered leakage. The reader may wonder what could happen if this assumption were dropped. First of all, we observe that it could make sense, i.e. in certain cases we could argue that the probabilistic knowledge associated to the secret choices (and its dependence on the observables) *could be considered as part of the leakage*. In the cases a and b of the cocaine auction example in Section 6, for instance, one may want to consider the information that we can deduce about the secrets (the identities of the bidder) from the observables

(the increments of the seller) as a leak due to the protocol. Our framework can encompass also this case, and the model remains the same. Nevertheless, the leakage would be represented by the mutual information rather than by the directed one.

In some other cases the flow of information from the observables to the secrets may even be considered as a consequence of the active attacks of an adversary, which uses the observables to modify the probability of the secrets. In this case the leakage would be divided in two parts: the one due to the protocol, represented by  $I(A^T \rightarrow B^T)$ , and the one due to the attacks of the adversaries, and represented by  $I(B^T \rightarrow A^T)$ . The total leakage would still be represented by the mutual information.

## 9 Future work

We would like to provide algorithms to compute the leakage and maximum leakage of interactive systems. These are rather challenging problems given the exponential growth of reaction functions (needed to compute the leakage) and the quantification over infinitely many reactors (given by the definition of maximum leakage in terms of capacity). One possible solution is to study the relation between deterministic schedulers and sequence of reaction functions. In particular, we believe that for each sequence of reaction functions and distribution over it there exists a probabilistic scheduler for the automata representation of the secret-nondeterministic IIHS. In this way, the problem of computing the leakage and maximum leakage would reduce to a standard probabilistic model checking problem (where the challenge is to compute probabilities ranging over infinitely many schedulers).

In addition, we plan to investigate measures of leakage for interactive systems other than mutual information and capacity.

We intend to study the applicability of our framework to the area of game theory. In particular, the interactive nature of games such as *Prisoner Dilemma* [21] and *Stag and Hunt* [23] (in their iterative versions) can be modeled as channels with memory and feedback following the techniques proposed in this work. Furthermore, (probabilistic) strategies can be encoded as reaction functions. In this way, optimal strategies are attained by reaction functions maximizing the leakage of the channel.

## Acknowledgement

We wish to thank the anonymous reviewers and Frank Valencia for their useful comments.

## References

1. M. S. Alvim, M. E. Andrés, and C. Palamidessi. Information flow in interactive systems. In P. Gastin and F. Laroussinie, editors, *Proceedings of the 21st International Conference on Concurrency Theory (CONCUR 2010)*, volume 6269 of *Lecture Notes in Computer Science*, pages 102–116. Springer, 2010.

2. M. E. Andrés, C. Palamidessi, P. van Rossum, and G. Smith. Computing the leakage of information-hiding systems. In J. Esparza and R. Majumdar, editors, *Proceedings of the Sixteenth International Conference on Tools and Algorithms for the Construction and Analysis of Systems (TACAS)*, volume 6015 of *Lecture Notes in Computer Science*, pages 373–389. Springer, 2010.
3. A. Bohannon, B. C. Pierce, V. Sjöberg, S. Weirich, and S. Zdancewic. Reactive noninterference. In E. Al-Shaer, S. Jha, and A. D. Keromytis, editors, *Proceedings of the 2009 ACM Conference on Computer and Communications Security, CCS 2009, Chicago, Illinois, USA, November 9-13, 2009*, pages 79–90. ACM, 2009.
4. K. Chatzikokolakis, C. Palamidessi, and P. Panangaden. Anonymity protocols as noisy channels. *Inf. and Comp.*, 206(2–4):378–401, 2008.
5. D. Clark, S. Hunt, and P. Malacaria. Quantified interference for a while language. In *Proceedings of the Second Workshop on Quantitative Aspects of Programming Languages (QAPL 2004)*, volume 112 of *Electronic Notes in Theoretical Computer Science*, pages 149–166. Elsevier Science B.V., 2005.
6. M. R. Clarkson, A. C. Myers, and F. B. Schneider. Belief in information flow. *Journal of Computer Security*, 17(5):655–701, 2009.
7. T. M. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley & Sons, Inc., second edition, 2006.
8. Y. Deng, T. Chothia, C. Palamidessi, and J. Pang. Metrics for action-labelled quantitative transition systems. In *Proceedings of the Third Workshop on Quantitative Aspects of Programming Languages (QAPL 2005)*, volume 153 of *Electronic Notes in Theoretical Computer Science*, pages 79–96. Elsevier Science Publishers, 2006. <http://www.lix.polytechnique.fr/~catuscia/papers/Metrics/QAPL/gts.pdf>.
9. J. Desharnais, R. Jagadeesan, V. Gupta, and P. Panangaden. The metric analogue of weak bisimulation for probabilistic processes. In *Proceedings of the 17th Annual IEEE Symposium on Logic in Computer Science*, pages 413–422. IEEE Computer Society, 2002.
10. Ebay website. <http://www.ebay.com/>.
11. Ebid website. <http://www.ebid.net/>.
12. R. G. Gallager. *Information Theory and Reliable Communication*. John Wiley & Sons, New York, NY, 1968.
13. J. W. Gray, III. Toward a mathematical foundation for information flow security. In *Proceedings of the 1991 IEEE Computer Society Symposium on Research in Security and Privacy (SSP '91)*, pages 21–35, Washington - Brussels - Tokyo, May 1991. IEEE.
14. L. Kantorovich. On the transfer of masses (in Russian). *Doklady Akademii Nauk*, 5(1):1–4, 1942. Translated in *Management Science*, 5(1):1–4, 1958.
15. B. Köpf and D. A. Basin. An information-theoretic model for adaptive side-channel attacks. In P. Ning, S. D. C. di Vimercati, and P. F. Syverson, editors, *Proceedings of the 2007 ACM Conference on Computer and Communications Security, CCS 2007, Alexandria, Virginia, USA, October 28-31, 2007*, pages 286–296. ACM, 2007.
16. P. Malacaria. Assessing security threats of looping constructs. In M. Hofmann and M. Felleisen, editors, *Proceedings of the 34th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2007, Nice, France, January 17-19, 2007*, pages 225–235. ACM, 2007.
17. J. L. Massey. Causality, feedback and directed information. In *Proceedings of the 1990 International Symposium on Information Theory and its Applications*, November 1990.
18. J. McLean. Security models and information flow. In *SSP'90*, pages 180–189. IEEE, 1990.
19. Mercadolibre website. <http://www.mercadolibre.com/>.
20. J. K. Millen. Hookup security for synchronous machines. In *Proceedings of the 3rd IEEE Computer Security Foundations Workshop (CSFW)*, pages 84–90, 1990.

21. W. Poundstone. *Prisoners Dilemma*. Doubleday NY, 1992.
22. R. Segala. *Modeling and Verification of Randomized Distributed Real-Time Systems*. PhD thesis, Massachusetts Institute of Technology, June 1995. Tech. Rep. MIT/LCS/TR-676.
23. B. Skyrms. *The Stag Hunt and the Evolution of Social Structure*. Cambridge University Press, 2003.
24. G. Smith. On the foundations of quantitative information flow. In L. de Alfaro, editor, *Proc. of the 12th Int. Conf. on Foundations of Software Science and Computation Structures*, volume 5504 of *LNCS*, pages 288–302, York, UK, 2009. Springer.
25. F. Stajano and R. J. Anderson. The cocaine auction protocol: On the power of anonymous broadcast. In *Information Hiding*, pages 434–447, 1999.
26. W. Stallings. *Data and Computer Communications*. Prentice Hall, eighth edition, 2006.
27. S. Subramanian. Design and verification of a secure electronic auction protocol. In *Proceedings of the 17th IEEE Symposium on Reliable Distributed Systems*, pages 204–210, Los Alamitos, CA, USA, 1998. IEEE Computer Society.
28. A. Tanenbaum. *Computer Networks*. Prentice Hall, second edition, 1989.
29. S. Tatikonda and S. K. Mitter. The capacity of channels with feedback. *IEEE Transactions on Information Theory*, 55(1):323–349, 2009.
30. F. van Breugel and J. Worrell. Towards quantitative verification of probabilistic transition systems. In F. Orejas, P. G. Spirakis, and J. van Leeuwen, editors, *Proceedings of the 28th International Colloquium on Automata, Languages and Programming (ICALP)*, volume 2076 of *Lecture Notes in Computer Science*, pages 421–432. Springer, 2001.
31. W. Vickrey. Counterspeculation, Auctions, and Competitive Sealed Tenders. *The Journal of Finance*, 16(1):8–37, 1961.

## Appendix

### A: An example illustrating the the encoder/decoder design

In this section we consider again the erasure channel of Example 2 to show how the enriched model of channels with memory and feedback can be used to transmit the message, and in particular how the feedback can be used to design the encoder. We assume that the set  $\mathcal{W}$  of possible messages consists of all finite sequences of bits. The role of the code functions is to encode the message  $W$  into a suitable representation for the stochastic kernels within the channel. The input alphabet for the stochastic kernels is  $\mathcal{A} = \{0, 1\}$  and the output alphabet is  $\mathcal{B} = \{0, 1, e\}$ , where the special output symbol  $e$  signals that a bit was erased. We assume that at most  $T$  uses of the channel are allowed and we use  $t$ , with  $1 \leq t \leq T$ , to represent the  $t^{\text{th}}$  time step.

We consider a sort of memory that depends only on the input history and we abstract from its specific form by defining a function  $\mu : \wp_f(\mathcal{A}^t) \mapsto [0, 1]$  that maps each possible input history to a correction factor to be added to (or subtracted from) a base probability value. We compute the contribution of  $\mu$  to the base values using arithmetics modulo 2, in such a way that the resulting values are still a probability distribution. More precisely, the stochastic kernels are defined as follows.

$$\begin{aligned}
p(\beta_t = 0 | \alpha^{t-1} 0, \beta^{t-1}) &= 0.8 - \mu(\alpha^{t-1}) \\
p(\beta_t = 1 | \alpha^{t-1} 0, \beta^{t-1}) &= 0 \\
p(\beta_t = e | \alpha^{t-1} 0, \beta^{t-1}) &= 0.2 + \mu(\alpha^{t-1}) \\
p(\beta_t = 0 | \alpha^{t-1} 1, \beta^{t-1}) &= 0 \\
p(\beta_t = 1 | \alpha^{t-1} 1, \beta^{t-1}) &= 0.8 - \mu(\alpha^{t-1}) \\
p(\beta_t = e | \alpha^{t-1} 1, \beta^{t-1}) &= 0.2 + \mu(\alpha^{t-1})
\end{aligned} \tag{14}$$

Correspondingly, the general form of the channel matrix for each time  $1 \leq t \leq T$  is shown in Table 10.

	0	1	e
$\alpha_t = 0, \beta^{t-1}$	$0.8 - \mu(\alpha^{t-1})$	0	$0.2 + \mu(\alpha^{t-1})$
$\alpha_t = 1, \beta^{t-1}$	0	$0.8 - \mu(\alpha^{t-1})$	$0.2 + \mu(\alpha^{t-1})$

**Table 10.** General form of channel matrix for  $1 \leq t \leq T$ .

The code functions are chosen at time  $t = 0$  based on the message to be transmitted. For illustration purposes, let us suppose that the message is the sequence of three bits  $W = 011$ . The other cases of  $W$  are analogous.

At time  $t = 1$ , the channel is used for its first time and the feedback history so far is empty  $\beta^0 = \epsilon$ . The encoder selects the input symbol  $\alpha_0 = 0$ , as in (15).

$$f_1[W = 011](\beta^0 = \epsilon) = 0 \tag{15}$$

At time  $t = 2$ , the feedback history consists of only one symbol, and in principle the possibilities are either  $\beta^1 = 0$ ,  $\beta^1 = 1$  or  $\beta^1 = e$ . In the first case, the first bit was successfully transmitted and the encoder can go on to the second bit of the message. By the way the channel is defined, the second case is not really possible, so it is not important how the reaction function is defined for this case. We will denote this indifference by attributing to the function the symbol  $x$  instead of a 0 or a 1. In the last case,  $\beta^1 = e$ , the first bit was erased and the encoder tries to retransmit the bit 0. We can write it formally as below.

$$\begin{aligned}
f_2[W = 011](\beta^1 = 0) &= 1 \\
f_2[W = 011](\beta^1 = 1) &= x \\
f_2[W = 011](\beta^1 = e) &= 0
\end{aligned} \tag{16}$$

At time  $t = 3$  the feedback histories allowed by the channel are  $\beta^2 \in \{01, 0e, e0, ee\}$  (the other ones have zero probability). In the first case,  $\beta^2 = 01$  the two first bits of the message have been transmitted correctly and the encoder can send the third bit. If  $\beta^2 = 0e$ , the transmission of the first bit was successful, but the second bit was erased and needs to be resent. In the case  $\beta^2 = e0$ , the first bit was erased in the first try but was successfully transmitted in the second try, so now the encoder can move to the

second bit of the message. In the last case,  $\beta^2 = ee$ , the two tries were unsuccessful and the encoder still needs to transmit the first bit of the message. Formally:

$$\begin{aligned}
f_3[W = 011](\beta^2 = 00) &= x \\
f_3[W = 011](\beta^2 = 01) &= 0 \\
f_3[W = 011](\beta^2 = 0e) &= 1 \\
f_3[W = 011](\beta^2 = 10) &= x \\
f_3[W = 011](\beta^2 = 11) &= x \\
f_3[W = 011](\beta^2 = 1e) &= x \\
f_3[W = 011](\beta^2 = e0) &= 1 \\
f_3[W = 011](\beta^2 = e1) &= x \\
f_3[W = 011](\beta^2 = ee) &= 0
\end{aligned} \tag{17}$$

We can easily extend the construction of code functions  $f_t$  for  $3 \leq t \leq T$  using this encoding scheme.

The decoder is very simple: once all time steps  $1, \dots, T$  have taken place, it just takes the output whole trace  $\beta^T$  and removes the occurrences of the erased bit symbol  $e$  in order to recover the original message.

Table 11 shows a concrete example of a possible behavior of binary erasure channel with memory and feedback in a scenario where the message is  $W = 011$  and the channel can be used  $T = 3$  times. Note that in this particular example the maximum uses of the channel is achieved before the whole message is successfully sent: the decoder can recover only the two first bits of the original message.

## B: Normalization of IIHS trees

In this section we will address the problem of *normalizing* an IIHS in such a way it is compatible with the assumptions made along the paper. The process of normalization described bellow is general enough to be applied to any IIHS without loss of generality or expressive power.

Consider a general IIHS  $\mathcal{J} = (M, \mathcal{A}, \mathcal{B})$  with  $M = (Q, \mathcal{L}, \hat{s}, \vartheta)$ , where  $\mathcal{L} = \mathcal{A} \cup \mathcal{B}$ . Assume that we are interested only in executions that involve up to  $T$  interactions, i.e  $T$  uses of the system, with one secret taking place and one observable produced at each time.

In the normalization process, we unfold the automaton up to level  $2T$ , since there is one secret symbol and one observable symbol for each step. We also extend the secret alphabet  $\mathcal{A}$  with a new symbol  $a_* \notin \mathcal{A}$  and the observable alphabet  $\mathcal{B}$  with a new symbol  $b_* \notin \mathcal{B}$ . These new symbols will be used as placeholders when we need to re-balance the tree. Let  $\mathcal{A}' = \mathcal{A} \cup \{a_*\}$  and  $\mathcal{B}' = \mathcal{B} \cup \{b_*\}$ .

For a given level  $t$  let  $\text{Labels}(\mathcal{J}, t)$  be the set of all labels of transitions that can be performed with a non-zero probability from the states at the  $t^{\text{th}}$  level of the automaton. Formally:

$$\text{Labels}(\mathcal{J}, t) \equiv \{\ell \in \mathcal{L} \mid \exists \sigma, s. |\sigma| = t, \text{last}(\sigma) \xrightarrow{\ell} s\}$$

The normalization of the IIHS  $\mathcal{J}$  leads to an equivalent IIHS  $\mathcal{J}' = (M', \mathcal{A}', \mathcal{B}')$ , where  $M' = (Q', \mathcal{L}', \hat{s}', \vartheta')$  and  $\mathcal{L}' = \mathcal{A}' \cup \mathcal{B}'$ ; and such that, for every  $1 \leq t \leq 2T$ :

Time $t$	Code functions $f_t(\beta^{t-1})$	Feedback history $\beta^{t-1}$	Encoder $\alpha_t = f_t[W](\beta^{t-1})$	Channel $p(\beta_t \alpha^t, \beta^{t-1})$	Decoder $\hat{W} = \gamma(\beta^T)$
$t = 0$	Code functions for $W = 011$ are selected.	————	————	————	————
$t = 1$	As in (15)	$\epsilon$	$\alpha_1 = f_1[W = 011](\epsilon)$ $= 0$	According to $p(\beta_1 0, \epsilon)$ produces $\beta_1 = \epsilon$	————
$t = 2$	As in s (16)	$\epsilon$	$\alpha_2 = f_2[W = 011](\epsilon)$ $= 0$	According to $p(\beta_2 00, \epsilon)$ produces $\beta_2 = 0$	————
$t = 3$	As in s (17)	$\epsilon 0$	$\alpha_3 = f_3[W = 011](\epsilon 0)$ $= 1$	According to $p(\beta_3 001, \epsilon 0)$ produces $\beta_3 = 1$	————
$t = 4$	————	————	————	————	Decoded message $\hat{W} = \gamma(\beta^3 = \epsilon 01)$ $= 01$

**Table 11.** A possible evolution of the binary channel with time, for  $W = 011$  and  $T = 3$

1.  $\text{Labels}(\mathcal{J}', t) \subseteq \mathcal{A}'$  or  $\text{Labels}(\mathcal{J}', t) \subseteq \mathcal{B}'$ ;
2.  $\text{Labels}(\mathcal{J}', t) \subseteq \mathcal{A}'$  iff  $\text{Labels}(\mathcal{J}', t + 1) \subseteq \mathcal{B}'$ , for  $1 \leq t \leq T - 1$ ;
3.  $\text{Labels}(\mathcal{J}', 1) \subseteq \mathcal{A}'$ ;

Condition 1 states that each level consists of either the secret actions only, or the observable actions only. Condition 2 states that secret and observable levels alternate. Condition 3 says that the automaton starts with a secret level.

The proof is straightforward. First, the new symbols  $a_*$  and  $b_*$  are placeholders for the absence of a secret and observable symbol, respectively. If in a given level  $t$  we want to have only secret symbols, we can postpone the occurrences of observable symbols at this level as follows: add  $a_*$  to the secret level and “move” all the observable symbols to the subtree of  $a_*$ . Figure 8 exemplifies the local transformations we need to make on the tree.

Note that in 8(b) the introduction of new nodes changed the probabilities of the transitions in the tree. In general, to normalize a secret level we need to introduce  $a_*$  in order to postpone the observable symbols, and the probabilities change as follows:

1. For every  $a_i$ ,  $1 \leq i \leq n$ , the associated probability is maintained as  $p'_{a_i} = p_{a_i}$ ;
2. The probability of the new symbol  $a_*$  is introduced as  $p_{a_*} = \sum_{i=0}^m p_{b_i}$ ;

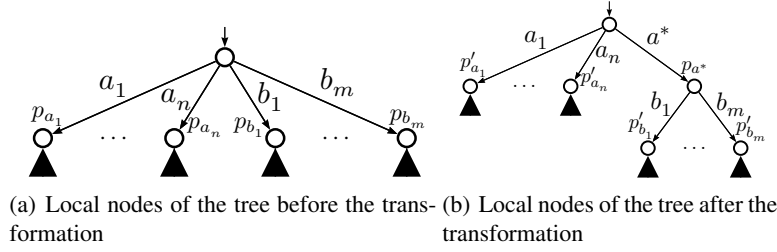


Fig. 8. Local transformation on an IIHS tree

3. If  $p_{a_*} \neq 0$ , then for  $1 \leq i \leq m$ , the associated probability of  $b_j$  is updated to  $p'_{b_j} = p_{b_j}/p_{a_*} = p_{b_j}/\sum_{k=0}^m p_{b_k}$ . If  $p_{a_*} = 0$ , then  $p'_{b_j} = 0$ , for  $1 \leq i \leq m$ , and  $p_{b_*} = 1$ .

The subtrees of each node of the original tree are preserved as they are, until we apply the same transformation to them. If a node does not have a subtree (i.e, no descendants), we create a subtree by adding all the possible actions in  $\mathcal{B}$  with probability 0, and the action  $b_*$  with probability 1.

If we are normalizing an observable level, the same rules apply, guarding the proper symmetry between secrets and observables. We then proceed on the same way on the deeper levels of the tree. Figure 9 shows an example of a full transformation on a tree (for the sake of readability, we omit the levels where only  $a_* = 1$  or  $b_* = 1$ ).

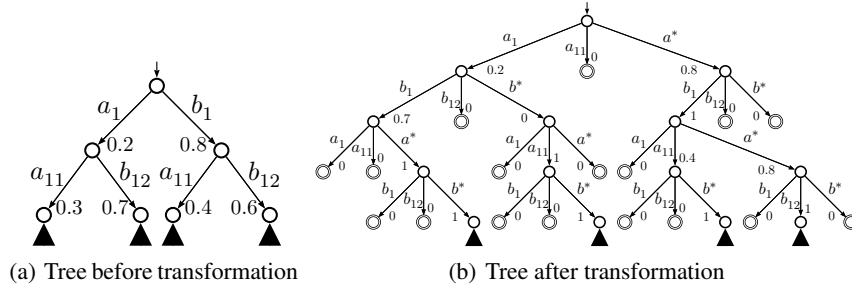


Fig. 9. Transformation on an IIHS tree

### C: Extended Cocaine Auction Protocol Example

We shall now extend the example of our approach applied to a real system. In Section 6 we introduced the Cocaine Auction Protocol and showed how to formalize one instance of it as an IIHS (Figure 7). We have also already defined the stochastic kernels for this example.

The next step is to construct all the possible reaction functions  $\{\varphi_t(\beta^{t-1})\}_{t=1}^T$ . As seen in Section 4.2, the reaction functions are the correspondent to the encoder in the



channel. They take the feedback story and decide how the world is going to react to this situation. Table 12 contains the reaction functions for each time  $t \leq 2$ .

(a) All 3 reaction functions  $\varphi_1$

$\beta^0$	$f_{1(1)}$	$f_{1(2)}$	$f_{1(3)}$
$\emptyset$	<i>Candlemaker</i>	<i>Scarface</i>	$a_*$

(b) All 27 reaction functions  $\varphi_2(\beta^1)$

$\beta^1$	$f_{2(1)}(\beta^1)$	$f_{2(2)}(\beta^1)$	$f_{2(3)}(\beta^1)$	$f_{2(4)}(\beta^1)$	$f_{2(5)}(\beta^1)$	$f_{2(6)}(\beta^1)$	$f_{2(7)}(\beta^1)$
$inc_1$	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Candlemaker</i>
$inc_2$	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Scarface</i>	<i>Scarface</i>	<i>Scarface</i>	$a_*$
$b_*$	<i>Candlemaker</i>	<i>Scarface</i>	$a_*$	<i>Candlemaker</i>	<i>Scarface</i>	$a_*$	<i>Candlemaker</i>
$\beta^1$	$f_{2(8)}(\beta^1)$	$f_{2(9)}(\beta^1)$	$f_{2(10)}(\beta^1)$	$f_{2(11)}(\beta^1)$	$f_{2(12)}(\beta^1)$	$f_{2(13)}(\beta^1)$	$f_{2(14)}(\beta^1)$
$inc_1$	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Scarface</i>	<i>Scarface</i>	<i>Scarface</i>	<i>Scarface</i>	<i>Scarface</i>
$inc_2$	$a_*$	$a_*$	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Scarface</i>	<i>Scarface</i>
$b_*$	<i>Scarface</i>	$a_*$	<i>Candlemaker</i>	<i>Scarface</i>	$a_*$	<i>Candlemaker</i>	<i>Scarface</i>
$\beta^1$	$f_{2(15)}(\beta^1)$	$f_{2(16)}(\beta^1)$	$f_{2(17)}(\beta^1)$	$f_{2(18)}(\beta^1)$	$f_{2(19)}(\beta^1)$	$f_{2(20)}(\beta^1)$	$f_{2(21)}(\beta^1)$
$inc_1$	<i>Scarface</i>	<i>Scarface</i>	<i>Scarface</i>	<i>Scarface</i>	$a_*$	$a_*$	$a_*$
$inc_2$	<i>Scarface</i>	$a_*$	$a_*$	$a_*$	<i>Candlemaker</i>	<i>Candlemaker</i>	<i>Candlemaker</i>
$b_*$	$a_*$	<i>Candlemaker</i>	<i>Scarface</i>	$a_*$	<i>Candlemaker</i>	<i>Scarface</i>	$a_*$
$\beta^1$	$f_{2(22)}(\beta^1)$	$f_{2(23)}(\beta^1)$	$f_{2(24)}(\beta^1)$	$f_{2(25)}(\beta^1)$	$f_{2(26)}(\beta^1)$	$f_{2(27)}(\beta^1)$	-
$inc_1$	$a_*$	$a_*$	$a_*$	$a_*$	$a_*$	$a_*$	-
$inc_2$	<i>Scarface</i>	<i>Scarface</i>	<i>Scarface</i>	$a_*$	$a_*$	$a_*$	-
$b_*$	<i>Candlemaker</i>	<i>Scarface</i>	$a_*$	<i>Candlemaker</i>	<i>Scarface</i>	$a_*$	-

**Table 12.** Reaction functions for the cocaine auction example.

Now we need to define the reactor, i.e., the probability distribution on reaction functions. Corollary 1 shows that we can do so by using the following equations:

$$p(\varphi_1) = p(\alpha_1 | \alpha^0, \beta^0) = p(\alpha_1)$$

$$p(\varphi_t | \varphi^{t-1}) = \prod_{\beta^{t-1}} p(\varphi_t(\beta^{t-1}) | \varphi^{t-1}(\beta^{t-2}), \beta^{t-1}), \quad 2 \leq t \leq T$$

For instance,  $p(f_{1(1)}) = p(\text{Candlemaker}) = p_1$ . In the same way,  $p(f_{1(2)}) = p(\text{Scarface}) = p_2$  and  $p(f_{1(3)}) = p(a_*) = p_3$ .

Let us take as an example the calculation of  $p(f_{2(6)} | f_{1(3)})$ :

$$\begin{aligned}
p(f_{2(6)}|f_{1(1)}) &= \prod_{\beta^1} p(f_{2(6)}(\beta^1)|\varphi_{1(1)}, \beta^1) \\
&= p(f_{2(6)}(inc_1)|Candlemaker, inc_1) \cdot p(f_{2(6)}(inc_2)|Candlemaker, inc_2) \\
&\quad p(f_{2(6)}(b_*)|Candlemaker, b_*) \\
&= p(Candlemaker|Candlemaker, inc_1) \cdot p(Scarface|Candlemaker, inc_2) \\
&\quad p(a_*|Candlemaker, b_*) \\
&= p_9 \cdot p_{13} \cdot 1 \\
&= p_9 p_{13} \tag{18}
\end{aligned}$$

Note that some reaction functions can have probability 0, which is consistent with the probabilistic automaton. For instance:

$$\begin{aligned}
p(f_{2(25)}|f_{1(3)}) &= \prod_{\beta^1} p(f_{2(25)}(\beta^1)|\varphi_{1(3)}, \beta^1) \\
&= p(f_{2(25)}(inc_1)|a_*, inc_1) \cdot p(f_{2(25)}(inc_2)|a_*, inc_2) \cdot p(f_{2(25)}(b_*)|a_*, b_*) \\
&= p(a_*|a_*, inc_1) \cdot p(a_*|a_*, inc_2) \cdot p(Candlemaker|a_*, b_*) \\
&= 1 \cdot 1 \cdot 0 \\
&= 0 \tag{19}
\end{aligned}$$