

Rate-Distortion Theory for Secrecy Systems

Curt Schieler and Paul Cuff

Abstract—Secrecy in communication systems is measured herein by the distortion that an adversary incurs. The transmitter and receiver share secret key, which they use to encrypt communication and ensure distortion at an adversary. A model is considered in which an adversary not only intercepts the communication from the transmitter to the receiver, but also potentially has side information. Specifically, the adversary may have causal or noncausal access to a signal that is correlated with the source sequence or the receiver’s reconstruction sequence. The main contribution is the characterization of the optimal tradeoff among communication rate, secret key rate, distortion at the adversary, and distortion at the legitimate receiver. It is demonstrated that causal side information at the adversary plays a pivotal role in this tradeoff. It is also shown that measures of secrecy based on normalized equivocation are a special case of the framework.

Index Terms—Rate-distortion theory, information-theoretic secrecy, shared secret key, causal disclosure, soft covering lemma, equivocation.

I. INTRODUCTION

In “Communication Theory of Secrecy Systems” [6], Shannon regarded a communication system as perfectly secret if the source and the eavesdropped message are statistically independent. The secrecy system studied in [6] is referred to as the “Shannon cipher system” and is depicted in Fig. 1. A necessary and sufficient condition for perfect secrecy is that the number of secret key bits per source symbol exceeds the entropy of the source. When the amount of key is insufficient, one must relax the requirement of statistical independence and invite new measures of secrecy.

One common way of measuring sub-perfect secrecy is with equivocation, the conditional entropy $H(X|M)$ of the source given the public message. The use of equivocation as a measure of secrecy was considered in the original work on the wiretap channel in [7] and [8], and it continues today. Although a distortion-based approach to secrecy might appear incomparable at first glance, it turns out that equivocation (when normalized by blocklength) becomes a special case of the framework developed here, under the proper choice of distortion measure.

In this work, we study an information-theoretic measure of secrecy that is directly inspired by rate-distortion theory. Whereas the objective in classical rate-distortion theory is to minimize a receiver’s distortion for a given rate of communication, our goal is to maximize an eavesdropper’s distortion for a given rate of secret key. If we relax the requirement of

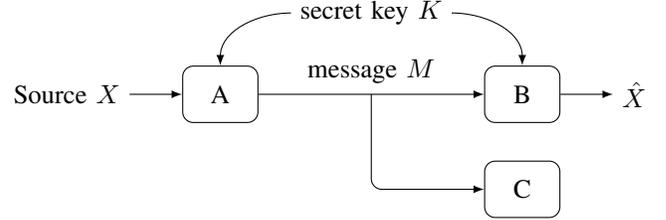


Fig. 1: The Shannon cipher system. Nodes A, B, and C are the transmitter, receiver, and eavesdropper, respectively.

lossless communication in Shannon’s cipher system, then our goal is to maximize an eavesdropper’s distortion for a given secret key rate, communication rate, and distortion tolerance at the receiver. Although there are a variety of secrecy systems other than Shannon’s cipher system (such as a wiretap channel [7] or distributed correlated sources [9], [10]), this paper is concerned exclusively with settings involving shared secret key, a single discrete memoryless source, and a noiseless channel. Moreover, we focus on block codes in the regime of blocklength tending to infinity.

When distortion is used as a measure of secrecy, we are implicitly viewing an eavesdropper in the same way that one views a receiver in a standard rate-distortion setting – as an active participant whose goal is to produce a sequence that is statistically correlated with the source sequence. Because he plays an active role, the eavesdropper is thought of as an adversarial entity. To ensure robustness, we will design the communication and encryption schemes against the worst-case adversarial strategy; that is, we wish to maximize the minimum distortion attainable by an adversary.

The study of information-theoretic secrecy via rate-distortion theory was initiated by Yamamoto in [11], in which the rate-distortion region was characterized for the special setting in which no secret key is available. Later, in “Rate-Distortion Theory for the Shannon Cipher System” [12], Yamamoto considered the exact problem we have heretofore described, but only obtained an inner and outer bound on the achievable rate-key-distortion region.¹ In this paper, we characterize the region; however, it is not our main focus. The following example serves to illustrate the care that should be exercised in a distortion-based approach to secrecy and motivates our primary investigation, which is centered around

¹The inner bound provided in [12] is precisely the region expressed in (49) of this work. Corollary 4 shows that this performance is achievable even if additional information is available to the eavesdropper, but it is suboptimal for the problem at hand. The outer bound in [12] makes use of two auxiliary variables, but with the appropriate selection can be shown to be equivalent to the trivial bound in (48), which in fact we show to be achievable. To show that the outer bound in [12] is trivial, the variables U and V can be selected as follows. Let U be independent of X and Y and uniformly distributed on $\{1, \dots, |\mathcal{X}|\}$. Let $V = U + X$ modulo $|\mathcal{X}|$.

This work was supported in part by the National Science Foundation under Grants CCF-1116013 and CCF-1017431, and also by the Air Force Office of Scientific Research under Grant FA9550-12-1-0196. Portions of this paper were presented in [1], [2], [3], [4], [5].

The authors are with the Department of Electrical Engineering, Princeton University, Princeton, NJ, 08544 USA (email: schieler@princeton.edu; cuff@princeton.edu).

a salient feature of our model referred to as *causal disclosure*.

A. One-bit secrecy and causal disclosure

Consider an n -bit i.i.d. source sequence $X^n \triangleq (X_1, \dots, X^n)$ with $X_i \sim \text{Bern}(1/2)$. Suppose common randomness $K \sim \text{Bern}(1/2)$ is available to the transmitter and receiver; that is, there is one bit of shared secret key. Now suppose the transmitter uses K to encrypt X^n by transmitting the n -bit message \tilde{X}^n , where $\tilde{X}_i = X_i \oplus K$. In other words, he flips all of the bits of X^n if $K = 1$, otherwise he simply sends X^n . Upon intercepting the public message \tilde{X}^n , the adversary produces a reconstruction Z^n and incurs expected distortion $\mathbb{E} \frac{1}{n} \sum_{i=1}^n d(X_i, Z_i)$, where $d(x, z)$ is a per-letter distortion measure. If $d(x, z) = \mathbf{1}\{x \neq z\}$, then an optimal strategy for the adversary is to simply set $Z^n = \tilde{X}^n$, yielding an expected distortion of $1/2$. Observe that $1/2$ is also the maximum possible expected distortion that we could ever force on the adversary, regardless of the amount of secret key available! It appears as though we have maximized secrecy by only using one bit of secret key for an arbitrarily long n -bit source. However, this view is severely misleading because the adversary actually knows a great deal about X^n , namely that it is one of only two candidate sequences.

This example demonstrates the potential fragility of using distortion to measure secrecy without recognizing the ramifications. For, although maximum secrecy (in the distortion sense) is attained, it vanishes altogether if the adversary views just one true bit of the source sequence (the bit allows him to determine whether or not to flip the \tilde{X}^n sequence). In general, the consequences of this example apply to the setting that Yamamoto considered in [12]. An arbitrarily small rate of secret key is enough to guarantee maximum distortion, but such secrecy is weak in the sense that even a small amount of additional knowledge (for example, observation of a few source symbols) is enough for the adversary to completely identify the source sequence.

The way that we strengthen a distortion-based approach to secrecy is through an assumption of causal disclosure, in which we design codes under the supposition that the adversary has noisy (or noiseless) access to the past behavior of the system. For example, in the one-bit secrecy example we might assume that the adversary produces the i th reconstruction symbol Z_i based not only on the public message M , but also on the past source symbols X^{i-1} . Incidentally, such a modification to the standard rate-distortion theory setting does not change the theory, though it has a dramatic effect in this secrecy setting. Regardless of whether or not an adversary actually has access to such information, designing our encryption under the assumption that he does leads to a much more robust notion of secrecy. In particular, it is resistant to disruptions in secrecy like those exhibited in the example. Despite the “pessimistic” nature of the causal disclosure assumption, we find that the optimal tradeoff between secret key and distortion in this regime is reasonable and not degenerate.

The assumption of causal disclosure is relevant not only for the sake of robustness, but also for its natural interpretations.

In [13], an alternative view of rate-distortion theory was introduced in which source and reconstruction sequences are regarded as sequences of actions in a distributed system. Communication is used to coordinate the receiver’s actions with the transmitter’s actions (which are given by nature). In this context, an adversary can be viewed as an active participant in the system who produces a sequence of actions. With this interpretation, it is not unrealistic to assume that the adversary could have causal access to the system behavior. Depending on where the adversary is intercepting communication, he might be able to view the past actions of the transmitter or receiver (or both) and produce his current action accordingly.

We find that optimal communication in this setting is not only fundamentally different than that of other source coding problems (often requiring a stochastic decoder), but in fact lends itself to a simple interpretation of injecting artificial memoryless noise into the adversary’s received signal.

B. Organization

The content of this paper is as follows. In Section II, we describe the problem setup. In Section III, we present a generalized version of the one-bit secrecy example in which there is no assumption of causal disclosure. In Section IV, we state our main result, Theorem 1, in which causal disclosure is a primary assumption. Theorem 1 describes the optimal relationship among the communication rate, secret key rate, and distortion at the legitimate receiver and adversary. Section IV also establishes a number of relevant corollaries to Theorem 1 and provides several concrete examples of the corresponding information-theoretic tradeoff regions. In Section V, we demonstrate how normalized equivocation arises as a special case of the causal disclosure framework. In Section VI, we give the achievability proof of Theorem 1. The proof uses a stochastic “likelihood encoder” that enables tractable analysis when combined with a “soft covering lemma”. Afterward, we discuss several important properties and implications of the optimal communication scheme used in the proof. Section VII provides the converse proof of Theorem 1. In Section VIII, we consider some settings with noncausal disclosure that are not subsumed by Theorem 1, but that can be proved similarly. Lastly, Section IX gives results for settings involving causal disclosure with delay greater than one.

II. PRELIMINARIES

The communication system model used throughout is shown in Fig. 2. The transmitting node, Node A, observes an i.i.d. source sequence $X^n \triangleq (X_1, \dots, X_n)$, where X_i is distributed according to P_X . Nodes A and B share a source of common randomness $K \in \{1, \dots, 2^{nR_0}\}$, referred to as secret key, that is uniformly distributed and independent of X^n . Based on the source block X^n and the secret key K , Node A transmits a message $M \in \{1, \dots, 2^{nR}\}$ that is received without loss by Nodes B and C. Once M is delivered, all three nodes sequentially produce actions: in the i th step, Nodes A, B and C produce X_i, Y_i , and Z_i , respectively. Note that Node A has no control over his actions; they are simply given by X^n . At the other end, Node B produces Y_i based on the pair (M, K)

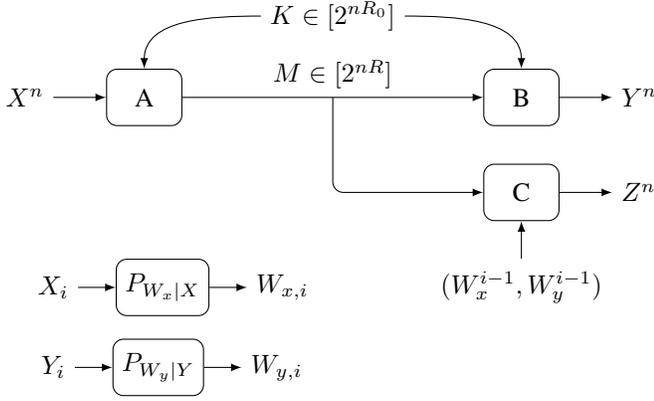


Fig. 2: Nodes A and B use secret key K and public communication M to coordinate against an adversarial Node C. At each step i , Node C can view the past behavior of the system, (W_x^{i-1}, W_y^{i-1}) , where W_x^n is the output of a memoryless channel $\prod P_{W_x|X}$ with input X^n , and W_y^n is the output of a memoryless channel $\prod P_{W_y|Y}$ with input Y^n .

and the adversarial Node C produces Z_i based on M and his observation of the past behavior of the system, (W_x^{i-1}, W_y^{i-1}) . At each step, the joint actions of the players incur a value $\pi(x, y, z)$, which represents symbol-wise payoff; the block-average payoff is given by

$$\frac{1}{n} \sum_{i=1}^n \pi(X_i, Y_i, Z_i). \quad (1)$$

Nodes A and B want to cooperatively maximize payoff, while Node C wants to minimize payoff through his actions Z^n . This payoff function can take the role of distortion incurred by Node C, corresponding to the secrecy metric described in the introduction. Note that instead of evaluating secrecy and coordination separately, which could be done with two payoff functions $\pi_1(x, y)$ and $\pi_2(x, y, z)$, we have unified them in a single function $\pi(x, y, z)$. Of course, the use of multiple payoff functions does have its own merits, and the results extend readily.

In Fig. 2, we depict noisy causal disclosure by (W_x^{i-1}, W_y^{i-1}) , where W_x^n is the output of a memoryless channel $\prod_{i=1}^n P_{W_x|X}$ with input X^n , and W_y^n is the output of a memoryless channel $\prod_{i=1}^n P_{W_y|Y}$ with input Y^n . Modeling the side information in this way covers a variety of scenarios. For example, if $P_{W_x|X}$ and $P_{W_y|Y}$ are identity channels, resulting in $(W_x, W_y) = (X, Y)$, then the adversary has full causal access (X^{i-1}, Y^{i-1}) . This is the strongest definition of secrecy in the causal disclosure framework and leads to the design of a thoroughly robust secrecy system. If $(W_x, W_y) = (\emptyset, \emptyset)$, then the adversary is completely blind to the past and only views the public message M ; this is the setting of [12], which does not include causal disclosure.

We remark that other strong security definitions involving side information leaks to the adversary can be found in [14], for example.

Throughout, we assume that the alphabets \mathcal{X} , \mathcal{Y} , and \mathcal{Z} are finite. We denote the set $\{1, \dots, m\}$ by $[m]$ and use $\Delta_{\mathcal{A}}$ to denote the probability simplex of distributions with alphabet

\mathcal{A} . The notation $X \perp Y$ indicates that the random variables X and Y are independent, and $X - Y - Z$ indicates a markov chain relationship.

Definition 1: An (n, R, R_0) code consists of an encoder $f : \mathcal{X}^n \times [2^{nR_0}] \rightarrow [2^{nR}]$ and a decoder $g : [2^{nR}] \times [2^{nR_0}] \rightarrow \mathcal{Y}^n$. More generally, we allow a stochastic encoder $P_{M|X^n, K}$ and a stochastic decoder $P_{Y^n|M, K}$. An (n, R, R_0) code is said to have blocklength n , communication rate R , and secret key rate R_0 .

Permitting stochastic decoders that use local randomization is crucial (in contrast to Wyner's wiretap channel, in which a stochastic *encoder* is needed). On the other hand, it is likely that the optimal encoder can be a deterministic function of the message and key, but this has not been shown. The proof of our main result uses a stochastic encoder and stochastic decoder.

Nodes A and B use an (n, R, R_0) code to coordinate against Node C. To ensure robustness, we consider the payoff that can be assured against the worst-case adversary, i.e., the max-min payoff. There are several ways to define the payoff criterion for a block, and we consider three: expected payoff, assured payoff, and symbol-wise minimum payoff. To distinguish among the three criteria, we use the monikers AVG, WHP, and MIN, respectively.

Definition 2: Fix a source distribution P_X , a symbol-wise payoff function $\pi : \mathcal{X} \times \mathcal{Y} \times \mathcal{Z} \rightarrow \mathbb{R}$, and causal disclosure channels $P_{W_x|X}$ and $P_{W_y|Y}$. For simplicity, denote the pair (W_x^n, W_y^n) by W^n . The triple (R, R_0, Π) is achievable if there exists a sequence of (n, R, R_0) codes such that

- Under the AVG criterion (expected payoff):

$$\liminf_{n \rightarrow \infty} \min_{\{P_{Z_i|M, W^{i-1}}\}_{i=1}^n} \mathbb{E} \frac{1}{n} \sum_{i=1}^n \pi(X_i, Y_i, Z_i) \geq \Pi. \quad (2)$$

- Under the WHP criterion (assured payoff):

$$\lim_{n \rightarrow \infty} \min_{\{P_{Z_i|M, W^{i-1}}\}_{i=1}^n} \mathbb{P} \left[\frac{1}{n} \sum_{i=1}^n \pi(X_i, Y_i, Z_i) \geq \Pi \right] = 1. \quad (3)$$

- Under the MIN criterion (symbol-wise minimum payoff):

$$\liminf_{n \rightarrow \infty} \min_{i \in [n]} \min_{P_{Z_i|M, W^{i-1}}} \mathbb{E} \pi(X_i, Y_i, Z) \geq \Pi. \quad (4)$$

Under the WHP criterion, the range of $\pi(x, y, z)$ is extended to include $-\infty$ so that lossless communication settings can be recovered.

Several remarks concerning the preceding definitions are in order.

- 1) Although WHP and MIN are incomparable, they are both stronger than AVG. However, it will be shown that all three criteria give rise to the same optimal tradeoff region.
- 2) In each of the criteria, we allow the adversary to employ his best set of probabilistic strategies $\{P_{Z_i|M, W^{i-1}}\}_{i=1}^n$ that minimize payoff. However, since expectation is linear in $P_{Z_i|M, W^{i-1}}$ for all i , the expectation is minimized by extreme points of the probability simplex; thus, we can assume that Node C uses a set of deterministic strategies, $\{z_i(m, w^{i-1})\}_{i=1}^n$.

- 3) It is assumed (although not explicit in the notation) that the adversary has full knowledge of the source distribution and the code that Nodes A and B use.
- 4) The optimal payoff does not increase if Node B is given direct causal access to Nodes A and C (i.e., if the decoder is given by $\{P_{Y_i|M,K,X^{i-1},Z^{i-1}}\}_{i=1}^n$ instead of simply $P_{Y^n|M,K}$). This is shown in Section VII in the converse proof of the main result.

Definition 3: The rate-payoff region \mathcal{R}_{AVG} is the closure of achievable triples (R, R_0, Π) under payoff criterion AVG. Regions \mathcal{R}_{WHP} and \mathcal{R}_{MIN} are defined in the same way.

III. ONE-BIT SECRECY, GENERALIZED

In this section, we expand on the scenario in which lossless communication is required between Nodes A and B and there is no causal disclosure of the system behavior to Node C. This is Yamamoto's setting in [12]. Although the result of this section is a special case of the main result in Theorem 1, it is an illustrative starting point.

For lossless communication, an additional achievability criterion is required, as stated below. Since X^n must equal Y^n with high probability, the payoff function is of the form $\pi(x, z)$. Thus, the achievability criteria for (R, R_0, Π) under the MIN payoff criterion (which is stronger than the AVG criterion) are that

$$\lim_{n \rightarrow \infty} \mathbb{P}[X^n \neq Y^n] = 0 \quad (5)$$

and

$$\liminf_{n \rightarrow \infty} \min_{i \in [n]} \min_{z(M)} \mathbb{E} \pi(X_i, z(M)) \geq \Pi. \quad (6)$$

Proposition 1: Fix P_X and $\pi(x, z)$. If lossless communication is required and there is no causal disclosure, then \mathcal{R}_{MIN} , the rate-payoff region under payoff criteria AVG and MIN is equal to

$$\left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq H(X) \\ R_0 \geq 0 \\ \Pi \leq \min_z \mathbb{E} \pi(X, z) \end{array} \right\}. \quad (7)$$

Thus, any positive rate of secret key² guarantees maximum secrecy (in the distortion sense), as Node C can achieve $\min_z \mathbb{E} \pi(X, z)$ by only knowing the source statistics. In fact, we now prove that each point in (7) can be achieved with key size $\mathcal{K} = [n]$ instead of $\mathcal{K} = [2^{nR_0}]$. This shows that even if the number of secret key bits is sublinear in the blocklength (in this case, $\log n$), one can still force the eavesdropper to incur the maximum distortion.³ As in the example of one-bit secrecy, such guarantees are shattered if even a small amount of source information is available to the adversary.

The following lemma is useful for the payoff analysis.

²Note that $R_0 = 0$ is only included in Proposition 1 because we defined the region as the *closure* of achievable triples. Furthermore, we remark that $R_0 = 0$ refers to a vanishing rate of secret key and is not the same as the absence of key.

³In [3], we show that an arbitrarily slow rate of increase is sufficient, even slower than $\log n$, under the AVG criterion.

Lemma 1: Let P_{XYZ} be a markov chain $X - Y - Z$, and f an arbitrary function. Then

$$\min_{g(x,y)} \mathbb{E} f(g(X, Y), Z) = \min_{g(y)} \mathbb{E} f(g(Y), Z). \quad (8)$$

Proof: We have

$$\begin{aligned} & \min_{g(x,y)} \mathbb{E} f(g(X, Y), Z) \\ &= \min_{g(x,y)} \sum_{x,y} P_{X,Y}(x,y) \mathbb{E}[f(g(X, Y), Z) | (X, Y) = (x, y)] \end{aligned} \quad (9)$$

$$= \sum_{x,y} P_{X,Y}(x,y) \min_g \mathbb{E}[f(g, Z) | (X, Y) = (x, y)] \quad (10)$$

$$\stackrel{(a)}{=} \sum_{x,y} P_{X,Y}(x,y) \min_g \mathbb{E}[f(g, Z) | Y = y] \quad (11)$$

$$= \min_{g(y)} \mathbb{E} f(g(Y), Z), \quad (12)$$

where (a) follows from the markovity assumption. \blacksquare

Now we prove Proposition 1.

Proof of Proposition 1: Converse. By the converse to the lossless source coding theorem, if (5) holds then we must have $R \geq H(X)$. To see that the payoff never exceeds $\min_z \mathbb{E} \pi(X, z)$, observe that the adversary can always let Z^n equal (z^*, \dots, z^*) , where

$$z^* = \operatorname{argmin}_z \mathbb{E} \pi(X, z). \quad (13)$$

Note that this converse argument holds for all three payoff criteria.

Achievability. Let $\varepsilon > 0$. Denote the empirical distribution (also referred to as the type) of a sequence x^n by P_{x^n} :

$$P_{x^n}(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{x_i = x\}. \quad (14)$$

The set of ε -typical sequences is defined as

$$\mathcal{T}_\varepsilon^n \triangleq \{x^n : |P_{x^n}(x) - P_X(x)| < \varepsilon P_X(x), \forall x \in \mathcal{X}\}. \quad (15)$$

To communicate, Nodes A and B use the set of ε -typical sequences as their codebook, just as in the standard proof of the lossless source coding theorem. If the source sequence X^n is typical, then the index of that codeword is the (pre-encrypted) message; if the source sequence is not typical, an arbitrary index is selected. Due to familiar properties of the size and probability of the typical set, the rate of communication is $(1 + \varepsilon)H(X)$ and the probability of error is

$$\mathbb{P}[X^n \neq Y^n] < \varepsilon \quad (16)$$

for large enough n .

The message will be encrypted using common randomness $K \sim \text{Unif}[n]$; this implies that the rate of secret key approaches zero as blocklength tends to infinity. In order to encrypt, we first partition $\mathcal{T}_\varepsilon^n$ into bins of size n (in a manner specified shortly), and use K to apply a one-time pad to the location of the source sequence X^n within the appropriate bin. More precisely, the encoder operates as follows: if X^n is ε -typical and is the L th sequence in the J th bin, then transmit the message $M = (J, L \oplus K)$, where \oplus indicates addition modulo n . By encrypting in this manner, the adversary knows

which bin X^n lies in (bin J), but does not know which of those n sequences it is, because L is independent of $L \oplus K$. Using the secret key, Node B can recover both J and L and produce the corresponding sequence.

The partitioning of $\mathcal{T}_\varepsilon^n$ is done according to the following equivalence relation:

$$x^n \sim y^n \text{ if } x^n \text{ is a cyclic permutation of } y^n. \quad (17)$$

Although the resulting partition can contain bins of size less than n , the number of such bins is small enough that we can ignore them without affecting the communication rate or (16). Thus, we assume that partitioning $\mathcal{T}_\varepsilon^n$ yields only bins of size n . Due to (17), it can be readily shown that each bin of size n has the following property.

Property 1: View the j th bin (denoted by b_j) as an $n \times n$ matrix whose columns are formed from the sequences in the bin. Then every row and column of the matrix has the same empirical distribution (denoted by P_j) and hence every row has the same probability (denoted by α_j) under the source distribution $\prod_{i=1}^n P_X(x_i)$.

This property is the crux of the proof; we offer the following intuition for why it implies that the eavesdropper suffers maximal distortion. The eavesdropper knows which bin X^n lies in, but does not know where it lies in the bin. Because of how we partitioned $\mathcal{T}_\varepsilon^n$, the eavesdropper's uncertainty is spread uniformly over the bin. To estimate X_i , the eavesdropper consults the i th row of the bin; however, Property 1 ensures that the empirical distribution of this row matches the type of the sequences in the bin, which in turn approximates the source distribution P_X (due to typicality). Therefore, the eavesdropper's estimate of X_i is based on no more than the original source statistics, which means that he suffers maximal distortion.

We now analyze the distortion precisely. For sufficiently large n , we have for all $i \in [n]$ that

$$\begin{aligned} & \min_{z^{(m)}} \mathbb{E} d(X_i, z(M)) \\ &= \min_{z^{(j,l)}} \mathbb{E} d(X_i, z(J, L \oplus K)) \end{aligned} \quad (18)$$

$$\stackrel{(a)}{=} \min_{z^{(j)}} \mathbb{E} d(X_i, z(J)) \quad (19)$$

$$= \min_{z^{(j)}} \sum_j \sum_{x^n \in b_j} p(x^n) d(x_i, z(j)) \quad (20)$$

$$\stackrel{(b)}{=} \sum_j \alpha_j \min_z \sum_{x^n \in b_j} d(x_i, z) \quad (21)$$

$$\stackrel{(c)}{=} \sum_j \alpha_j \min_z \sum_{x \in \mathcal{X}} n P_j(x) d(x, z) \quad (22)$$

$$\stackrel{(d)}{\geq} \sum_j n \alpha_j \min_z \sum_{x \in \mathcal{X}} (1 - \varepsilon) P_X(x) d(x, z) \quad (23)$$

$$= \mathbb{P}[X^n \in \mathcal{T}_\varepsilon^n] (1 - \varepsilon) \min_z \mathbb{E} d(X, z) \quad (24)$$

$$\geq (1 - \varepsilon)^2 \min_z \mathbb{E} d(X, z), \quad (25)$$

where (a) is due to $(X_i, J) \perp (L \oplus K)$ and Lemma 1, (b) and (c) are due to Property 1, and (d) follows from the definition of $\mathcal{T}_\varepsilon^n$. Thus, we have (6). ■

Discussion

Suppose Nodes A and B use the binning scheme just described in the proof of Proposition 1 to achieve maximum secrecy. What if, instead of eavesdropping only the public message, the adversary is also able to view the past behavior of the system, namely X^{i-1} ? Because of the structure of each bin (i.e., Property 1), knowledge of just the first symbol, $X_1 = x_1$, is enough for the adversary to narrow down the size of the list of candidate source sequences from n to approximately $n P_X(x_1)$. One can see that the adversary will be able to determine the true sequence quickly, well before the end of the block. In this manner, the adversary can take advantage of the causal disclosure to force the payoff to take on its minimum value instead of its maximum value. In general, causal disclosure benefits an adversary and gives rise to a nontrivial tradeoff between secret key and payoff. We remark that one of the key elements in the proof of the main result is that the benefits of causal disclosure can be voided if the right amount of secret key is available. In fact, it will become evident in Section VI that using secret key to sterilize the causal disclosure gives rise to the optimal tradeoff of secret key and payoff.

IV. MAIN RESULT

Our main result is the following.

Theorem 1: Fix P_X , $\pi(x, y, z)$, and causal disclosure channels $P_{W_x|X}$ and $P_{W_y|Y}$. Then \mathcal{R}_{AVG} , the closure of achievable (R, R_0, Π) under payoff criterion AVG, is equal to

$$\bigcup_{W_x - X - (U, V) - Y - W_y} \left\{ (R, R_0, \Pi) : \begin{aligned} & R \geq I(X; U, V) \\ & R_0 \geq I(W_x W_y; V|U) \\ & \Pi \leq \min_{z(u)} \mathbb{E} \pi(X, Y, z(U)) \end{aligned} \right\}, \quad (26)$$

where $|U| \leq |\mathcal{X}| + 2$ and $|V| \leq |\mathcal{X}||\mathcal{Y}|(|\mathcal{X}| + 2) + 1$. Furthermore,

$$\mathcal{R}_{\text{AVG}} = \mathcal{R}_{\text{WHP}} = \mathcal{R}_{\text{MIN}}. \quad (27)$$

We remark that the convexity of \mathcal{R}_{AVG} and \mathcal{R}_{WHP} can be shown from Definitions 2 and 3 by using a standard time-sharing argument. By (27), \mathcal{R}_{MIN} is also a convex set.

We now elaborate on several corollaries to Theorem 1 that are obtained through different choices of the causal disclosure channels $P_{W_x|X}$ and $P_{W_y|Y}$. To begin, we consider scenarios in which lossless communication is required between Nodes A and B.

A. Lossless communication

In the following, we require X^n to equal Y^n with high probability. That is, we introduce into Definition 2 the additional constraint

$$\lim_{n \rightarrow \infty} \mathbb{P}[X^n \neq Y^n] = 0. \quad (28)$$

Conveniently, (28) can be ensured by considering payoff criterion WHP with a payoff function $\pi(x, y, z)$ that evaluates to $-\infty$ when $x \neq y$.

Corollary 1: Fix P_X , $\pi(x, z)$, and causal disclosure channel $P_{W_x|X}$. If lossless communication is required (i.e., (28) is imposed), then the rate-payoff region \mathcal{R}_{WHP} is equal to

$$\bigcup_{U-X-W_x} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq H(X) \\ R_0 \geq I(W_x; X|U) \\ \Pi \leq \min_{z(u)} \mathbb{E} \pi(X, z(U)) \end{array} \right\}. \quad (29)$$

Proof: Define a payoff function

$$\bar{\pi}(x, y, z) \triangleq \begin{cases} \pi(x, z) & \text{if } x = y \\ -\infty & \text{if } x \neq y. \end{cases} \quad (30)$$

When $\Pi > -\infty$, it is easily verified that

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[\frac{1}{n} \sum_{i=1}^n \bar{\pi}(X_i, Y_i, Z_i) \geq \Pi - \varepsilon \right] = 1 \quad (31)$$

if and only if both of the following hold:

$$\lim_{n \rightarrow \infty} \mathbb{P}[X^n = Y^n] = 1 \quad (32)$$

$$\lim_{n \rightarrow \infty} \mathbb{P} \left[\frac{1}{n} \sum_{i=1}^n \pi(X_i, Z_i) \geq \Pi - \varepsilon \right] = 1. \quad (33)$$

Thus, \mathcal{R}_{WHP} (the region we seek the characterize) is obtained by invoking Theorem 1 with $W_y = \emptyset$. However, we want to simplify the region further. Denoting the region in (29) by \mathcal{S} , we now show that $\mathcal{R}_{\text{WHP}} = \mathcal{S}$.

Note that when $\Pi > -\infty$, we have

$$-\infty < \Pi \leq \min_{z(u)} \mathbb{E} \bar{\pi}(X, Y, z(U)), \quad (34)$$

which implies $X = Y$. When combined with the markov chain $X-(U, V)-Y$, this gives $H(X|UV) = 0$. Therefore, $\mathcal{R}_{\text{WHP}} \subseteq \mathcal{S}$ follows from writing

$$\begin{aligned} R &\geq I(X; U, V) = H(X) \\ R_0 &\geq I(W_x; V|U) = I(W_x; X, V|U) \geq I(W_x; X|U). \end{aligned}$$

To see $\mathcal{S} \subseteq \mathcal{R}_{\text{WHP}}$, let $V = Y = X$. \blacksquare

Corollary 1, in turn, spawns two important results. By invoking Corollary 1 with $W_x = \emptyset$, we recover Proposition 1 under WHP.

Corollary 2: Fix P_X and $\pi(x, z)$. If lossless communication is required and there is no causal disclosure, then the rate-payoff region \mathcal{R}_{WHP} is equal to

$$\left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq H(X) \\ R_0 \geq 0 \\ \Pi \leq \min_z \mathbb{E} \pi(X, z) \end{array} \right\}. \quad (35)$$

If we instead consider the disclosure channel $W_x = X$, we have the following.

Corollary 3: Fix P_X and $\pi(x, z)$. If lossless communication is required and X^{i-1} is disclosed, then the rate-payoff

region \mathcal{R}_{WHP} is equal to

$$\bigcup_{P_{U|X}} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq H(X) \\ R_0 \geq H(X|U) \\ \Pi \leq \min_{z(u)} \mathbb{E} \pi(X, z(U)) \end{array} \right\}. \quad (36)$$

B. Lossless communication example

In this section, we present a concrete example of the region in Corollary 3 (causal disclosure of Node A) and compare it to the region in Corollary 2 (no causal disclosure).

We first show that (36) can be written as a linear program. Since the constraint on R is fixed by the source distribution, we focus our attention on the boundary of the (R_0, Π) tradeoff, namely

$$\Pi(R_0) \triangleq \max_{\substack{P_{U|X}: \\ H(X|U) \geq R_0}} \min_{z(u)} \mathbb{E} \pi(X, z(U)). \quad (37)$$

Notice that this can be rewritten as

$$\Pi(R_0) = \max_{\substack{P_U, P_{X|U}: \\ \sum_u P_U P_{X|U} = P_X \\ H(X|U) \geq R_0}} \sum_u P_U(u) \min_z \mathbb{E}[\pi(X, z)|U = u]. \quad (38)$$

If we are able to restrict the set $\{P_{X|U=u}\}_{u \in \mathcal{U}}$ in the maximization to a finite set $\mathcal{P} \subseteq \Delta_{\mathcal{X}}$, then $\Pi(R_0)$ can be expressed as a linear program. Indeed, viewing the distribution P_U as a vector $p \in \mathbb{R}^{|\mathcal{P}|}$, (38) becomes

$$\begin{aligned} &\text{maximize} && d^\top p \\ &\text{subject to} && p \geq 0 \\ &&& \mathbf{1}^\top p = 1 \\ &&& T p = P_X \\ &&& h^\top p \leq R_0 \end{aligned} \quad (39)$$

where

- $T \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{P}|}$ is the transition matrix whose columns are the elements of \mathcal{P} .
- The vector $d \in \mathbb{R}^{|\mathcal{P}|}$ has entries

$$d_u = \min_z \mathbb{E}[\pi(X, z)|U = u], \quad u \in \mathcal{U}. \quad (40)$$

- The vector $h \in \mathbb{R}^{|\mathcal{P}|}$ has entries

$$h_u = H(X|U = u), \quad u \in \mathcal{U}. \quad (41)$$

To see why there is always a choice of finite \mathcal{P} such that the rate-payoff boundary is unaffected, consider the function $d : \Delta_{\mathcal{X}} \rightarrow \mathbb{R}$ defined by

$$d(p) = \min_z \mathbb{E}[\pi(X, z)], \text{ where } X \sim p. \quad (42)$$

Observe that $d(\cdot)$ is the boundary of a convex polytope because it is the minimum of $|\mathcal{Z}|$ linear functions (and \mathcal{Z} is finite). Define the set

$$\mathcal{P} = \{p \in \Delta_{\mathcal{X}} : d(p) \text{ is an extreme point of } d\} \quad (43)$$

Given a set of distributions $\{P_{X|U=u}\}_{u \in \mathcal{U}}$ that optimize (38), we can write each element $P_{X|U=u}$ as a convex combination of the distributions in \mathcal{P} while maintaining the value of

the objective. Furthermore, due to the concavity of the entropy function, the constraint on R_0 is still satisfied. Thus, \mathcal{P} is sufficient for the optimization.

In the particular case that the payoff function is hamming distance (i.e., $\pi(x, z) = \mathbf{1}\{x \neq z\}$), the set \mathcal{P} has a particularly convenient form:

$$\mathcal{P} = \{p \in \Delta_{\mathcal{X}} : p = \text{Unif}(\mathcal{A}) \text{ for some } \mathcal{A} \subseteq \mathcal{X}\}. \quad (44)$$

This allows us to give the following simple analytical expression for $\Pi(R_0)$. The proof is given in Appendix A.

Theorem 2: Fix P_X and let $\pi(x, z) = \mathbf{1}\{x \neq z\}$. Define the function $\phi(\cdot)$ as the linear interpolation of the points $(\log n, \frac{n-1}{n})$, $n \in \mathbb{N}$.⁴ Also, define

$$\pi_{\max} = 1 - \max_x P_X(x). \quad (45)$$

Then, the boundary of the rate-payoff region when lossless communication is required and X^{i-1} is disclosed can be written as

$$\Pi(R_0) = \min\{\phi(R_0), \pi_{\max}\}. \quad (46)$$

In Fig. 3, we illustrate Theorem 2 for an arbitrary source distribution. Note that when there is *no* causal disclosure and $\pi(x, z)$ is hamming distance, the payoff is given by Corollary 2 as

$$\min_z \mathbb{E} \pi(X, z) = 1 - \max_x P_X(x) = \pi_{\max}, \quad (47)$$

regardless of the rate of secret key. Comparing (47) with $\min\{\phi(R_0), \pi_{\max}\}$ demonstrates the effect of causal disclosure (see Fig. 3). In particular, we see that the assumption that the adversary does not view any of the true source bits can lead to a rather fragile guarantee of maximum secrecy. Indeed, at low rates of secret key, the gap that results from revealing the source causally is the difference between maximum secrecy and zero secrecy. This reduction in payoff is the price that is paid for increased robustness against an adversary (e.g., preventing pitfalls like those that we saw in the example of one-bit secrecy).

From Theorem 2, we also readily see that the payoff can saturate when $R_0 < H(X)$, which shows that maximum payoff is not the same as Shannon's perfect secrecy. For example, if $P_X = \{1/4, 1/4, 1/2\}$, then the maximum payoff of $1/2$ occurs at $R_0 = 1$, but $H(X) = 1.5$.

C. Lossy communication

In the previous section, the communication rate lay above $H(X)$ and did not affect the (R_0, Π) tradeoff. However, when the requirement of lossless communication is relaxed, all three quantities interact. There are four natural special cases that are obtained by setting W_x equal to \emptyset or X and setting W_y equal to \emptyset or Y . We denote the corresponding rate-payoff regions as \mathcal{R}_\emptyset , \mathcal{R}_A , \mathcal{R}_B , and \mathcal{R}_{AB} to distinguish which nodes' actions are causally revealed.

Corollary 4: Fix P_X and $\pi(x, y, z)$. In each of the following, the region holds under all three payoff criteria.

⁴Here n does not refer to blocklength.

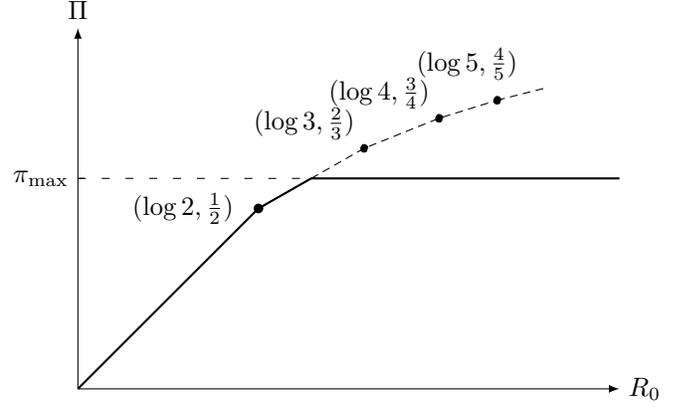


Fig. 3: Illustration of Theorem 2 for a generic source P_X with $1 - \max_x P_X(x) = \pi_{\max}$. The solid curve, $\Pi(R_0) = \min\{\phi(R_0), \pi_{\max}\}$, is the tradeoff between rate of secret key and payoff under the assumption of causal disclosure (Corollary 3). The loosely dashed line is π_{\max} , which also corresponds to the payoff when there is no causal disclosure (Corollary 2). The densely dashed curve is $\phi(R_0)$.

If there is no causal disclosure, then the rate-payoff region, \mathcal{R}_\emptyset , is equal to

$$\bigcup_{P_{Y|X}} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq I(X; Y) \\ R_0 \geq 0 \\ \Pi \leq \min_z \mathbb{E} \pi(X, Y, z) \end{array} \right\}. \quad (48)$$

If X^{i-1} is disclosed, then the rate-payoff region, \mathcal{R}_A , is equal to

$$\bigcup_{P_{Y,U|X}} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq I(X; Y, U) \\ R_0 \geq I(X; Y|U) \\ \Pi \leq \min_{z(u)} \mathbb{E} \pi(X, Y, z(u)) \end{array} \right\}. \quad (49)$$

If Y^{i-1} is disclosed, then \mathcal{R}_B is given by directly substituting $W_x = \emptyset$ and $W_y = Y$ in (26). Similarly, if (X^{i-1}, Y^{i-1}) is disclosed, then \mathcal{R}_{AB} is given by directly substituting $W_x = X$ and $W_y = Y$ in (26).

Proof: Setting $(W_x, W_y) = (\emptyset, \emptyset)$ in Theorem 1 gives \mathcal{R}_\emptyset . Denote the region in (48) by \mathcal{S} . If $(R, R_0, \Pi) \in \mathcal{R}_\emptyset$, then

$$R \geq I(X; U, V) = I(X; U, V, Y) \geq I(X; Y) \quad (50)$$

$$\Pi \leq \min_{z(u)} \mathbb{E} \pi(X, Y, z(U)) \leq \min_z \mathbb{E} \pi(X, Y, z), \quad (51)$$

which gives $\mathcal{R}_\emptyset \subseteq \mathcal{S}$. To see $\mathcal{S} \subseteq \mathcal{R}_\emptyset$, let $U = \emptyset$ and $V = Y$.

Setting $(W_x, W_y) = (X, \emptyset)$ in Theorem 1 gives \mathcal{R}_A . Denote the region in (49) by \mathcal{T} . If $(R, R_0, \Pi) \in \mathcal{R}_A$, then

$$R \geq I(X; U, V) = I(X; U, V, Y) \geq I(X; U, Y) \quad (52)$$

$$R_0 \geq I(X; V|U) = I(X; V, Y|U) \geq I(X; Y|U), \quad (53)$$

which gives $\mathcal{R}_A \subseteq \mathcal{T}$. To see $\mathcal{T} \subseteq \mathcal{R}_A$, let $V = Y$. ■

D. Lossy communication examples

In this section, we investigate concrete examples of Corollary 4 by considering the payoff function

$$\pi(x, y, z) = \mathbf{1}\{x = y, x \neq z\}. \quad (54)$$

For this choice, the block-average payoff is the fraction of symbols in a block that Nodes A and B are able to agree on and keep hidden from Node C.

We now present achievable regions for the cases of Corollary 4 when $P_X \sim \text{Bern}(1/2)$ and $\pi(x, y, z)$ is given by (54). The region that we give for R_\emptyset is optimal, and numerical computation suggests that the other regions are optimal as well. Setting $P_{Y|X} = \text{BSC}(\alpha)$, we have

$$\mathcal{R}_\emptyset = \bigcup_{\alpha \in [0, \frac{1}{2}]} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq 1 - h(\alpha) \\ R_0 \geq 0 \\ \Pi \leq \frac{1}{2}(1 - \alpha) \end{array} \right\}. \quad (55)$$

If we let $U = \emptyset$ and $P_{Y|X} = \text{BSC}(\alpha)$, then we have

$$\mathcal{R}_A \supseteq \bigcup_{\alpha \in [0, \frac{1}{2}]} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq 1 - h(\alpha) \\ R_0 \geq 1 - h(\alpha) \\ \Pi \leq \frac{1}{2}(1 - \alpha) \end{array} \right\}. \quad (56)$$

Letting $U = \emptyset$, $P_{Y|X} = \text{BSC}(\alpha)$, and $P_{V|Y} = \text{BSC}(\beta)$ gives

$$\mathcal{R}_B \supseteq \bigcup_{\alpha, \beta \in [0, \frac{1}{2}]} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq 1 - h(\alpha) \\ R_0 \geq 1 - h(\beta) \\ \Pi \leq \frac{1}{2}(1 - \alpha \star \beta) \end{array} \right\} \quad (57)$$

and also

$$\mathcal{R}_{AB} \supseteq \text{conv} \left(\bigcup_{\alpha, \beta \in [0, \frac{1}{2}]} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq 1 - h(\alpha) \\ R_0 \geq 1 + h(\alpha \star \beta) \\ \quad - h(\alpha) - h(\beta) \\ \Pi \leq \frac{1}{2}(1 - \alpha \star \beta) \end{array} \right\} \right). \quad (58)$$

where $\alpha \star \beta = \alpha(1 - \beta) + \beta(1 - \alpha)$ and $\text{conv}(\cdot)$ denotes the convex hull operation. Regions (56) and (57) are convex as given.

Several observations concerning the regions in Fig. 4 are in order. First, the minimum payoff is $1/4$, which occurs when there is no communication or secret key. This is achieved if Node B generates an i.i.d. sequence according to $\text{Bern}(1/2)$ and Node C produces an arbitrary sequence. The maximum payoff that can be guaranteed is $1/2$, because Node C can correctly guess X with probability one-half without any information. Second, note the strict containment from top to bottom: causal access to Node A (Fig. 4b) is better for the adversary than access to Node B (Fig. 4c), and the combination (Fig. 4d) is strictly better for him than Node A alone. Finally, observe the effect of having a higher secret key rate than communication rate, and vice versa. When Node A is causally revealed, the payoff is a function of $\min(R, R_0)$ and there is

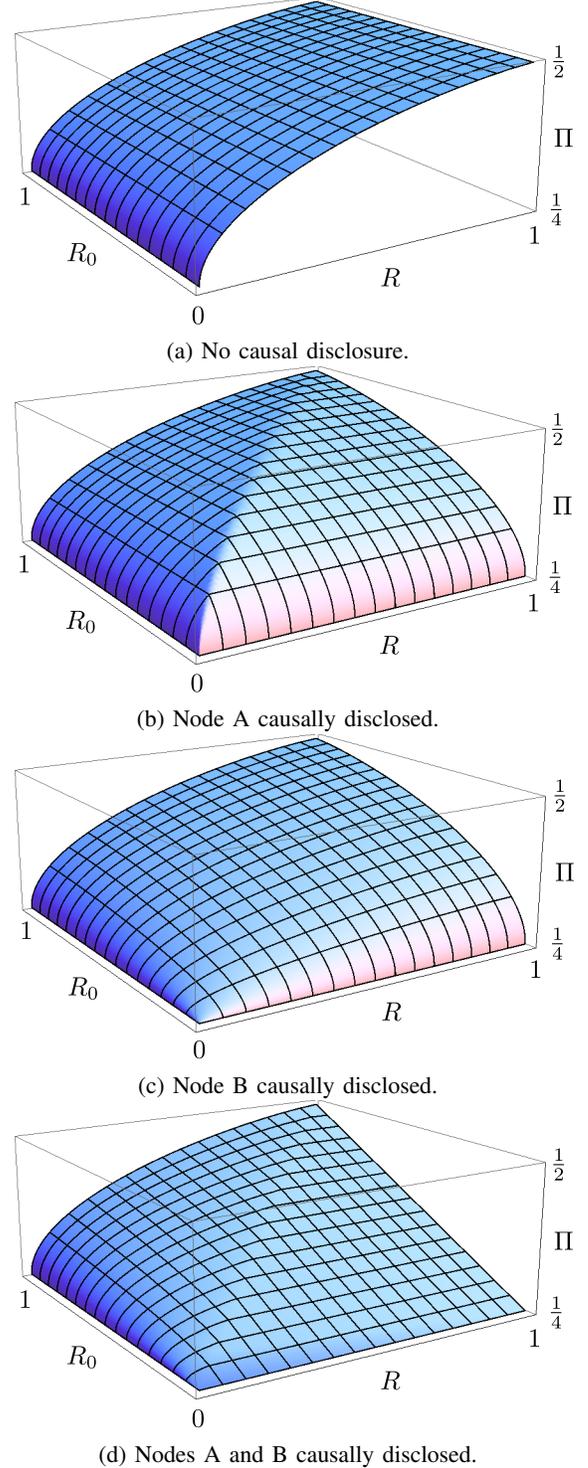


Fig. 4: Achievable regions of Corollary 4 for $P_X \sim \text{Bern}(1/2)$ and $\pi(x, y, z) = \mathbf{1}\{x = y, x \neq z\}$. Numerical computation suggests that these regions are optimal.

no advantage in having excess of either rate. However, when Node B is revealed, both $R_0 > R$ and $R > R_0$ result in higher payoff than $R = R_0$. When both nodes are revealed, an excess of secret key rate increases payoff.⁵ This phenomenon is particularly surprising because it means that secret key is useful even beyond the application of a one-time-pad to the communication.

V. EQUIVOCATION

In this section, we show that (normalized) equivocation-based measures of secrecy become a special case of the causal disclosure framework if we choose the payoff function to be a log-loss function. Relating distortion to conditional entropy via a log-loss function was done recently in the context of certain multiterminal source coding problems [15].

First, we remark that Theorem 1 can be readily extended to include multiple distortion functions. For example, if we wanted to separately evaluate coordination and secrecy, we could use two payoff functions $\pi_1(x, y)$ and $\pi_2(x, y, z)$. In this setting, it might be more natural to refer to distortion functions than payoff functions, with the goal of minimizing the distortion between Nodes A and B while maximizing the distortion between Nodes (A,B) and Node C. Then, the rate-distortion region becomes

$$\bigcup_{W_x - X - (U, V) - Y - W_y} \left\{ \begin{array}{l} (R, R_0, D_1, D_2) : \\ R \geq I(X; U, V) \\ R_0 \geq I(W_x W_y; V|U) \\ D_1 \geq \mathbb{E} d_1(X, Y) \\ D_2 \leq \min_{z(u)} \mathbb{E} d_2(X, Y, z(U)) \end{array} \right\}. \quad (59)$$

Now consider $(W_x, W_y) = (X, \emptyset)$ and a distortion function $d_2 : \mathcal{X} \times \mathcal{Y} \times \Delta_{\mathcal{X}} \rightarrow \mathbb{R}$ defined by

$$d_2(x, y, z) = \log \frac{1}{z(x)} \quad (60)$$

where z is a probability distribution on \mathcal{X} , and $z(x)$ denotes the probability of $x \in \mathcal{X}$ according to $z \in \Delta_{\mathcal{X}}$. With this choice, the distortion in criterion AVG can be written as

$$\begin{aligned} & \min_{\{P_{Z_i|M, X^{i-1}}\}_{i=1}^n} \mathbb{E} \frac{1}{n} \sum_{i=1}^n d_2(X_i, Y_i, Z_i) \\ &= \frac{1}{n} \sum_{i=1}^n \min_{P_{Z|M, X^{i-1}}} \mathbb{E} d_2(X_i, Y_i, Z) \end{aligned} \quad (61)$$

$$= \frac{1}{n} \sum_{i=1}^n \min_{P_{Z|M, X^{i-1}}} \mathbb{E} \log \frac{1}{Z(X_i)} \quad (62)$$

$$\stackrel{(a)}{=} \frac{1}{n} \sum_{i=1}^n H(X_i|M, X^{i-1}) \quad (63)$$

$$= \frac{1}{n} H(X^n|M), \quad (64)$$

where (a) is due to the Lemma 2 (given below). Thus, for the log-loss distortion function in (60), expected adversarial distortion under an assumption of causal disclosure simply becomes normalized equivocation.

⁵These relationships are not known to be true in general.

Lemma 2: Fix a pair of random variables (X, Y) and let $\mathcal{Z} = \Delta_{\mathcal{X}}$. Then

$$H(X|Y) = \min_{z: X-Y-Z} \mathbb{E} \log \frac{1}{Z(X)} \quad (65)$$

where $z(x)$ is the probability of x according to z .

Proof: If $X - Y - Z$, then

$$\mathbb{E} \log \frac{1}{Z(X)} \quad (66)$$

$$= \mathbb{E} \log \frac{1}{P_{X|Y}(X|Y)} + \mathbb{E} \log \frac{P_{X|Y}(X|Y)}{Z(X)} \quad (67)$$

$$= H(X|Y) + \sum_{y,z} P_{YZ}(y, z) D(P_{X|Y=y} || z) \quad (68)$$

$$\geq H(X|Y), \quad (69)$$

with equality if $z = P_{X|Y=y}$ for all (y, z) . ■

So far, we have focused on the equivocation of X^n ; however, one might be interested in $\frac{1}{n} H(Y^n|M)$ or $\frac{1}{n} H(X^n, Y^n|M)$, instead. In these cases, the rate-distortion-equivocation regions can again be recovered from Theorem 1 (via the form in (59)) by considering $(W_x, W_y) = (\emptyset, Y)$, $\mathcal{Z} = \Delta_{\mathcal{Y}}$ and

$$d_2(x, y, z) = \log \frac{1}{z(y)} \quad (70)$$

or $(W_x, W_y) = (X, Y)$, $\mathcal{Z} = \Delta_{\mathcal{X} \times \mathcal{Y}}$ and

$$d_2(x, y, z) = \log \frac{1}{z(x, y)}. \quad (71)$$

In all three cases, the regions can be simplified (in particular, the auxiliary random variable V can be eliminated). The results are given in the following theorem, part 1 of which was given by Yamamoto in [12].

Corollary 5: Fix P_X and $d(x, y)$. Let \mathcal{R} denote the closure of achievable pairs (R, R_0, D, E) .

1) If the equivocation criterion is

$$\liminf_{n \rightarrow \infty} \frac{1}{n} H(X^n|M) \geq E, \quad (72)$$

then

$$\mathcal{R} = \bigcup_{P_{Y|X}} \left\{ \begin{array}{l} (R, R_0, D, E) : \\ R \geq I(X; Y) \\ D \geq \mathbb{E} d(X, Y) \\ E \leq H(X) - [I(X; Y) - R_0]_+ \end{array} \right\}, \quad (73)$$

where $[x]_+ = \max\{0, x\}$.

2) If the equivocation criterion is

$$\liminf_{n \rightarrow \infty} \frac{1}{n} H(Y^n|M) \geq E, \quad (74)$$

then

$$\mathcal{R} = \bigcup_{X-U-Y} \left\{ \begin{array}{l} (R, R_0, D, E) : \\ R \geq I(X; U) \\ D \geq \mathbb{E} d(X, Y) \\ E \leq H(Y) - [I(Y; U) - R_0]_+ \end{array} \right\}. \quad (75)$$

3) If the equivocation criterion is

$$\liminf_{n \rightarrow \infty} \frac{1}{n} H(X^n, Y^n|M) \geq E, \quad (76)$$

then

$$\mathcal{R} = \bigcup_{X-U-Y} \left\{ \begin{array}{l} (R, R_0, D, E) : \\ R \geq I(X; U) \\ D \geq \mathbb{E} d(X, Y) \\ E \leq H(X, Y) - [I(X, Y; U) - R_0]_+ \end{array} \right\}. \quad (77)$$

Proof: We only prove part 2, as parts 1 and 3 follow similar arguments. First, fix $d_2(x, y, z)$ according to (70). Then, by Lemma 2,

$$\min_{z(u)} \mathbb{E} d_2(X, Y, z(U)) = H(Y|U). \quad (78)$$

From the discussion above, it is clear that \mathcal{R} is characterized by setting $(W_x, W_y) = (\emptyset, Y)$ in (59), yielding

$$\mathcal{R} = \bigcup_{X-(U,V)-Y} \left\{ \begin{array}{l} (R, R_0, D, E) : \\ R \geq I(X; U, V) \\ R_0 \geq I(Y; V|U) \\ D \geq \mathbb{E} d(X, Y) \\ E \leq H(Y|U) \end{array} \right\}. \quad (79)$$

Denote the region in (75) by \mathcal{S} . To see $\mathcal{R} \subseteq \mathcal{S}$, first consider $(R, R_0, D, E) \in \mathcal{R}$. Defining $U' \triangleq (U, V)$, we have

$$R \geq I(X; U, V) = I(X; U') \quad (80)$$

$$E \leq H(Y|U) = H(Y|U, V) + I(Y; V|U) \quad (81)$$

$$\leq H(Y|U') + R_0 \quad (82)$$

$$E \leq H(Y), \quad (83)$$

which implies $(R, R_0, D, E) \in \mathcal{S}$. To see $\mathcal{S} \subseteq \mathcal{R}$, let $(R, R_0, D, E) \in \mathcal{S}$. Define $V' \triangleq U$ and find a random variable U' such that $U' - U - (X, Y)$ form a markov chain and

$$H(Y|U') = H(Y) - [I(Y; U) - R_0]_+ \quad (84)$$

This is always possible because the right-hand side of (84) lies in the interval $[H(Y|U), H(Y)]$. Then, we can write

$$R \geq I(X; U) = I(X; U', V') \quad (85)$$

$$R_0 \geq H(Y|U') - H(Y|U) = I(Y; V'|U') \quad (86)$$

$$E \leq H(Y|U'), \quad (87)$$

which implies $(R, R_0, D, E) \in \mathcal{R}$. Thus, $\mathcal{R} = \mathcal{S}$. ■

VI. ACHIEVABILITY PROOF

A. Soft covering lemma

The primary tool used in the achievability proof of Theorem 1 is a so-called ‘‘soft covering lemma’’, a known result concerning the approximation of the output distribution of a channel.⁶ Various forms of the lemma have appeared in [17] and [18] and related notions from the perspective of random binning can be found in [19]. Several generalizations of the lemma (including a one-shot version) can be found in [16].

In brief, the most basic version of the soft covering lemma is as follows. Fix a joint distribution $P_{X,U}$. First, generate a random codebook of 2^{nR} independent codewords, each

drawn according to $\prod_{i=1}^n P_U(u_i)$. Select a codeword, uniformly at random, as the input to a memoryless channel $\prod_{i=1}^n P_{X|U}(x_i|u_i)$. The lemma states that if $R > I(X; U)$, then the distribution of the channel output X^n converges to $\prod_{i=1}^n P_X(x_i)$ in expected total variation distance, where the expectation is with respect to the random codebook.

A generalization of the soft covering lemma, presented shortly, will prove essential to the payoff analysis. Once we define a code by pairing a random codebook with a particular stochastic encoder and decoder, the soft covering lemma can be used to approximate the joint statistics of the system (i.e., the joint distribution on (X^n, M, K, Y^n, W^n) that is induced by the code) by an ‘‘idealized’’ distribution that has desirable properties. Having a tractable approximation of the joint distribution of the system is important because an adversary’s optimal strategy is dictated by a posterior distribution. For example, if an adversary tries to estimate the i th source symbol X_i based on his observations of causal disclosure X^{i-1} and the public message M , his optimal strategy is entirely determined by the posterior distribution of X_i given (X^{i-1}, M) . The approximating distribution that the soft covering lemma guarantees will provide a clear understanding of that posterior distribution and lead to a manageable payoff analysis.

Although the distribution approximation in the soft covering lemma holds for normalized and unnormalized divergence, we use the total variation version found in [16] and [18] because of the following properties that total variation enjoys.

Given two probability measures P and Q with common alphabet \mathcal{X} , the total variation distance between P and Q is defined by

$$\|P - Q\| = \sup_{A \in \mathcal{F}} |P(A) - Q(A)|, \quad (88)$$

where \mathcal{F} is the sigma algebra of the common alphabet.

Property 2: Total variation distance satisfies the following.

(a) If the support of P and Q is a countable set \mathcal{X} , then

$$\|P - Q\| = \frac{1}{2} \sum_{x \in \mathcal{X}} |P(\{x\}) - Q(\{x\})|. \quad (89)$$

(b) Let $\varepsilon > 0$ and let $f(x)$ be a function with bounded range of width $b > 0$. Then

$$\|P - Q\| < \varepsilon \implies |\mathbb{E}_P f(X) - \mathbb{E}_Q f(X)| < \varepsilon b, \quad (90)$$

where \mathbb{E}_P indicates that the expectation is taken with respect to the distribution P .

(c) For any P, Q , and Φ ,

$$\|P - Q\| \leq \|P - \Phi\| + \|\Phi - Q\|. \quad (91)$$

(d) Let $P_X P_{Y|X}$ and $Q_X P_{Y|X}$ be two joint distributions with common channel $P_{Y|X}$. Then

$$\|P_X P_{Y|X} - Q_X P_{Y|X}\| = \|P_X - Q_X\|. \quad (92)$$

(e) Let P_X and Q_X be marginal distributions of P_{XY} and Q_{XY} . Then

$$\|P_X - Q_X\| \leq \|P_{XY} - Q_{XY}\|. \quad (93)$$

We require the following generalization of the soft covering lemma.

⁶The name ‘‘soft covering lemma’’ was given in [16]. The same lemma has also been referred to as the ‘‘resolvability lemma’’ and ‘‘cloud-mixing lemma’’.

Definition 4: Let $\{P_{X^n, Y^n}\}_{n=1}^{\infty}$ be a sequence of joint distributions. The sup-information rate of this sequence is defined as

$$\bar{I}(X; Y) \triangleq \limsup_{n \rightarrow \infty} \frac{1}{n} i_{P_{X^n, Y^n}}(X^n; Y^n), \quad (94)$$

where

$$\limsup_{n \rightarrow \infty} W_n \triangleq \inf\{\tau : \mathbb{P}[W_n > \tau] \rightarrow 0\} \quad (95)$$

and

$$i_{P_{X, Y}}(a; b) \triangleq \log \frac{P_{X, Y}(a, b)}{P_X(a)P_Y(b)}. \quad (96)$$

The function $i_{P_{X, Y}}(a; b)$ is called the information density.

Lemma 3 ([16, Corollary VII.4], [18]): Let $\{P_{X^n, Y^n}\}_{n=1}^{\infty}$ be a sequence of joint distributions. Let $\mathcal{C}^{(n)}$ be a random codebook of 2^{nR} sequences in \mathcal{X}^n , each drawn independently according to P_{X^n} and indexed by $m \in [2^{nR}]$. Let Q_{Y^n} denote the output distribution of the channel when the input is selected from $\mathcal{C}^{(n)}$ uniformly at random; that is,

$$Q_{Y^n}(y^n) = 2^{-nR} \sum_{m \in [2^{nR}]} P_{Y^n|X^n}(y^n|X^n(m)). \quad (97)$$

If $R > \bar{I}(X; Y)$, then

$$\lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{C}^{(n)}} \|Q_{Y^n} - P_{Y^n}\| = 0, \quad (98)$$

where $\mathbb{E}_{\mathcal{C}^{(n)}}$ indicates that the expectation is with respect to the random codebook.⁷ Furthermore, the convergence in (98) occurs exponentially quickly with n if the distribution $P_{X^n Y^n}$ is memoryless.

We now begin the achievability proof of Theorem 1 by specifying the random codebook, stochastic encoder, and stochastic decoder.

B. Design of codebook, encoder, and decoder

In the statement of Theorem 1, we are given disclosure channels $P_{W_x|X}$ and $P_{W_y|Y}$. For simplicity, we treat the channels as a single channel⁸ defined by

$$P_{W|XY} \triangleq P_{W_x|X} P_{W_y|Y}. \quad (99)$$

Thus, we denote (W_x^n, W_y^n) by W^n and the causal disclosure by W^{i-1} . The memoryless channel from (X^n, Y^n) to W^n is denoted by

$$P_{W^n|X^n Y^n} \triangleq \prod_{i=1}^n P_{W|XY}. \quad (100)$$

Given a source distribution P_X and a disclosure channel $P_{W|XY}$, fix a distribution

$$P_{XUVY} = P_X P_{UV|X} P_{Y|UV} P_{W|XY}. \quad (101)$$

⁷Because the codebook is random, the output distribution Q_{Y^n} is a random variable taking values on $\Delta_{\mathcal{Y}^n}$. One way to notate this is through the use of conditional distributions (i.e., write $Q_{Y^n|\mathcal{C}^{(n)}}$), but we choose to suppress such notation in order to simplify the presentation.

⁸The decomposition of the channel $P_{W|XY}$ into two channels does not play a role in the achievability proof. The reason Theorem 1 does not feature a generic channel $P_{W|XY}$ is that a matching converse proof has not been supplied.

Note that this distribution satisfies the markov chain $X - (U, V) - Y$. Fix a communication rate $R > I(X; U, V)$ and a secret key rate $R_0 > I(W; V|U)$.

Random codebook: Generate a random superposition codebook in the following manner. First, generate a codebook $\mathcal{C}_U^{(n)}$ of 2^{nR} codewords from \mathcal{U}^n i.i.d. according to $\prod_{i=1}^n P_U$. These codewords are indexed by $m \in [2^{nR}]$. Then, for each codeword $U^n(m) \in \mathcal{C}_U^{(n)}$, generate a codebook $\mathcal{C}_V^{(n)}(m)$ of 2^{nR_0} codewords from \mathcal{V}^n i.i.d. according to $\prod_{i=1}^n P_{V|U=U_i(m)}$. These codewords are indexed by $(m, k), k \in [2^{nR_0}]$. Thus, we have

$$\mathcal{C}_U^{(n)} = (U^n(1), \dots, U^n(m), \dots, U^n(2^{nR})) \quad (102)$$

and

$$\mathcal{C}_V^{(n)}(m) = (V^n(m, 1), \dots, V^n(m, k), \dots, V^n(m, 2^{nR_0})). \quad (103)$$

We refer to the entire superposition codebook as $\mathcal{C}^{(n)}$.

Likelihood encoder: For a fixed superposition codebook, the encoder is a stochastic likelihood encoder defined by

$$P_{M|X^n K}(m|x^n, k) \propto \prod_{i=1}^n P_{X|U, V}(x_i|u_i(m), v_i(m, k)), \quad (104)$$

where \propto indicates that an appropriate normalization factor is required to make $P_{M|X^n K}$ a valid conditional probability distribution. Eqn. (104) says that the probability of (x^n, k) being mapped to the index m is proportional to the probability that x^n is the output of the memoryless “test channel” $P_{X|UV}$ with input $(u^n(m), v^n(m, k))$. The reason for this choice of encoder will become clear shortly.

Decoder: The decoder is stochastic and is defined by

$$P_{Y^n|MK}(y^n|m, k) \triangleq \prod_{i=1}^n P_{Y|UV}(y_i|u_i(m), v_i(m, k)). \quad (105)$$

The random codebook, likelihood encoder, and decoder comprise the code and induce a joint distribution on the system that is given by

$$P_{X^n MK Y^n W^n} = P_{X^n} P_K P_{M|X^n K} P_{Y^n|MK} P_{W^n|X^n Y^n}, \quad (106)$$

where P_{X^n} is i.i.d. according to P_X , and P_K is uniform over $[2^{nR_0}]$.

C. The approximating distribution Q and its property

We now use the soft covering lemma (Lemma 3) to yield an approximation to the system-induced distribution $P_{X^n MK Y^n W^n}$. The idealized distribution that we are concerned with is described by Fig. 5 and defined explicitly as

$$Q_{X^n MK Y^n W^n} \triangleq Q_{X^n MK} P_{Y^n|MK} P_{W^n|X^n Y^n}, \quad (107)$$

where $Q_{X^n MK}$ is given by

$$Q(x^n, m, k) \triangleq 2^{n(R+R_0)} \prod_{i=1}^n P_{X|UV}(x_i|U_i(m), V_i(m, k)). \quad (108)$$

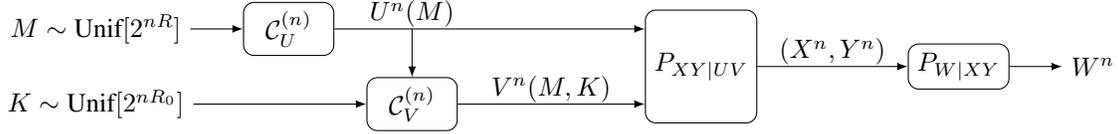


Fig. 5: Process that defines $Q_{X^n M K Y^n W^n}$. The pair (M, K) indexes a pair of codewords $(U^n(M), V^n(M, K))$ in the superposition random codebook. The codeword pair is passed through a memoryless channel $P_{XY|UV} = P_{X|UV}P_{Y|UV}$ to get (X^n, Y^n) . Then (X^n, Y^n) is passed through a memoryless channel $P_{W|XY}$ to get W^n .

Observe that the definitions of $P_{Y^n|MK}$ and $P_{W^n|X^n Y^n}$, combined with the factorization of P_{XUVY^W} in (101), allow us to write Q as

$$Q(x^n, m, k, y^n, w^n) = 2^{-n(R+R_0)} \prod_{i=1}^n P_{XYW|UV}(x_i, y_i, w_i | U_i(m), V_i(m, k)), \quad (109)$$

which corresponds to the process depicted in Fig. 5.

The stochastic likelihood encoder was defined intentionally so that $Q_{M|X^n K} = P_{M|X^n K}$. In fact, the only difference between P and Q lies in the marginal distribution of (X^n, K) . Indeed, notice that we can write

$$Q_{X^n M K Y^n W^n} \triangleq Q_{X^n M K} P_{Y^n|MK} P_{W^n|X^n Y^n} \quad (110)$$

$$= Q_{X^n K} P_{M|X^n K} P_{Y^n|MK} P_{W^n|X^n Y^n} \quad (111)$$

$$= Q_{X^n K} P_{M Y^n W^n | X^n K}. \quad (112)$$

Therefore, we can show that $P_{X^n M K Y^n W^n} \approx Q_{X^n M K Y^n W^n}$ by demonstrating that $P_{X^n K} \approx Q_{X^n K}$. This is accomplished using the soft covering lemma.

Lemma 4: If $R > I(X; U, V)$, then

$$\lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{C}^{(n)}} \left\| P_{X^n M K Y^n W^n} - Q_{X^n M K Y^n W^n} \right\| = 0. \quad (113)$$

Proof of Lemma 4:

$$\mathbb{E}_{\mathcal{C}^{(n)}} \left\| P_{X^n M K Y^n W^n} - Q_{X^n M K Y^n W^n} \right\| \stackrel{(a)}{=} \mathbb{E}_{\mathcal{C}^{(n)}} \left\| P_{X^n K} - Q_{X^n K} \right\| \quad (114)$$

$$= \mathbb{E}_{\mathcal{C}^{(n)}} \left\| P_{X^n} P_K - Q_{X^n|K} P_K \right\| \quad (115)$$

$$\stackrel{(b)}{=} 2^{-nR_0} \sum_{k=1}^{2^{nR_0}} \mathbb{E}_{\mathcal{C}^{(n)}} \left\| P_{X^n} - Q_{X^n|K=k} \right\| \quad (116)$$

$$= \mathbb{E}_{\mathcal{C}^{(n)}} \left\| P_{X^n} - Q_{X^n|K=1} \right\| \quad (117)$$

$$\stackrel{(c)}{=} o(1). \quad (118)$$

The justification for the steps is as follows:

- (a) Eqn. (112) and Property 2d of total variation.
- (b) Property 2a of total variation.
- (c) $R > I(X; U, V)$ and the soft covering lemma (Lemma 3). Notice that P_{X^n} is i.i.d. according to P_X and $Q_{X^n|K=1}$ is the output distribution of the memoryless channel $P_{X|UV}$ acting on a (sub)codebook of size 2^{nR} .

Approximating P by Q will allow us to analyze the payoff as if Q governs the joint statistics of the system. If the rate of secret key is large enough, the structure of Q will allow us to

argue that the causal disclosure W^{i-1} is actually useless to the eavesdropper and that his best strategy for estimating (X_i, Y_i) is based solely on $U^n(M)$. The crucial property of Q that enables this argument is given in the following lemma. The proof, which relies on the soft covering lemma, is relegated to Appendix B.

Lemma 5: If $R_0 > I(W; V|U)$, there exists $\alpha \in (0, 1]$ such that

$$\lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{C}^{(n)}} \left\| Q_{M W^n X_B Y_B} - \widehat{Q}_{M W^n X_B Y_B} \right\| = 0, \quad (119)$$

where

$$\widehat{Q}_{M W^n X_B Y_B} \triangleq Q_M \cdot \left(\prod_{i=1}^n P_{W|U=U_i(M)} \right) \cdot \left(\prod_{i \in \mathcal{B}} P_{XY|W, U=U_i(M)} \right) \quad (120)$$

and \mathcal{B} is any subset of $[n]$ of size $|\mathcal{B}| \leq \lfloor \alpha n \rfloor$.

To see the significance of \widehat{Q} , first consider $\mathcal{B} = \emptyset$ and $W = (X, Y)$, so that

$$\widehat{Q}_{M X^n Y^n}(m, x^n, y^n) = 2^{-nR} \prod_{i=1}^n P_{XY|U}(x_i, y_i | U_i(m)). \quad (121)$$

Recall that $W = (X, Y)$ implies direct causal disclosure of Nodes A and B; that is, the adversary has access to (M, X^{i-1}, Y^{i-1}) at step i . From (121), we see that $\widehat{Q}_{X^n Y^n|M}$ is a memoryless channel from the codeword $U^n(M)$ to the pair (X^n, Y^n) . In particular, this implies

$$(X_i, Y_i) - U_i(M) - (M, X^{i-1}, Y^{i-1}), \forall i \in [n]. \quad (122)$$

Therefore, the adversary's best estimate of (X_i, Y_i) only depends on $U_i(M)$ and is not improved by the causal disclosure. We have essentially created an artificial noisy channel from the intercepted codeword $U^n(M)$ to the pair (X^n, Y^n) , a property which not only greatly simplifies the payoff analysis, but is interesting independent of the causal disclosure problem. We discuss this effect and some of its implications after completing the achievability proof.

For general W , consider $\mathcal{B} = \{i\}$. In this case, Lemma 5 demonstrates that Q approximately satisfies the markov chain

$$(X_i, Y_i) - U_i(M) - (M, W^{i-1}), \quad (123)$$

and again we see that adversary's estimate of (X_i, Y_i) only depends on $U_i(M)$ and is not improved by the causal disclosure. However, it turns out that the property in (123) is not quite strong enough for the analysis of the WHP payoff criterion, which is why Lemma 5 is concerned with sub-blocks $(X_{\mathcal{B}}, Y_{\mathcal{B}})$ of size linearly increasing with n .

D. Analysis of the MIN payoff criterion

We first combine Lemmas 4 and 5 to demonstrate the existence of a codebook that ensures certain distribution approximations hold simultaneously for all $i \in [n]$.

Lemma 6: There exists a sequence of codebooks such that

$$\lim_{n \rightarrow \infty} \max_{i \in [n]} \|P_{MW^n X_i Y_i} - \widehat{Q}_{MW^n X_i Y_i}\| = 0 \quad (124)$$

and

$$\lim_{n \rightarrow \infty} \max_{i \in [n]} \|\widehat{Q}_{u_i(M)} - P_U\| = 0. \quad (125)$$

Proof of Lemma 6: First, for all $i \in [n]$ we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{C}^{(n)}} \|P_{MW^n X_i Y_i} - \widehat{Q}_{MW^n X_i Y_i}\| \\ & \stackrel{(a)}{\leq} \mathbb{E}_{\mathcal{C}^{(n)}} \|P_{MW^n X_i Y_i} - Q_{MW^n X_i Y_i}\| \\ & \quad + \mathbb{E}_{\mathcal{C}^{(n)}} \|Q_{MW^n X_i Y_i} - \widehat{Q}_{MW^n X_i Y_i}\| \end{aligned} \quad (126)$$

$$\begin{aligned} & \stackrel{(b)}{\leq} \mathbb{E}_{\mathcal{C}^{(n)}} \|P_{X^n M K Y^n W^n} - Q_{X^n M K Y^n W^n}\| \\ & \quad + \mathbb{E}_{\mathcal{C}^{(n)}} \|Q_{MW^n X_i Y_i} - \widehat{Q}_{MW^n X_i Y_i}\| \end{aligned} \quad (127)$$

$$\stackrel{(c)}{=} O(e^{-\gamma n}) \quad (128)$$

for some $\gamma > 0$. Steps (a) and (b) use Properties 2c and 2e of total variation distance, respectively. Step (c) follows from Lemmas 4 and 5, and the fact that the convergence in the soft covering lemma occurs exponentially quickly with n .

Next, we invoke Lemma 3 to show that, for all $i \in [n]$,

$$\mathbb{E}_{\mathcal{C}^{(n)}} \|\widehat{Q}_{U_i(M)} - P_U\| = O(e^{-\beta n}), \quad (129)$$

for some $\beta > 0$. The soft covering lemma applies because:

- $\widehat{Q}_{U_i(M)}$ is the output distribution of the identity channel acting on a “codebook” of 2^{nR} “codewords” generated i.i.d. according to P_U – the “codebook” consists of $(U_i(1), \dots, U_i(2^{nR}))$. Furthermore P_U is the output distribution when the input distribution is P_U , because the channel is the identity channel.
- The rate requirement is trivially satisfied because $R > 0$ and

$$\limsup_{n \rightarrow \infty} \frac{1}{n} i_{P_U}(U_i; U_i) \leq \lim_{n \rightarrow \infty} \frac{1}{n} \log |\mathcal{U}| = 0. \quad (130)$$

Combining (128) and (129), we can write

$$\begin{aligned} & \lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{C}^{(n)}} \left[\sum_{i=1}^n \|P_{X^n Y_i M} - \widehat{Q}_{X^n Y_i M}\| \right. \\ & \quad \left. + \sum_{i=1}^n \|\widehat{Q}_{U_i(M)} - P_U\| \right] = 0. \end{aligned} \quad (131)$$

It is straightforward to verify that this fact implies the statement of the lemma. \blacksquare

With Lemma 6 in hand, we proceed with the analysis of the MIN payoff criterion. Let $\Pi \leq \min_{z(u)} \mathbb{E} \pi(X, Y, z(U))$. For

all $i \in [n]$, we have

$$\begin{aligned} & \min_{z(m, w^{i-1})} \mathbb{E}_P \pi(X_i, Y_i, z(M, W^{i-1})) \\ & \stackrel{(a)}{=} \min_{z(m, w^{i-1})} \mathbb{E}_{\widehat{Q}} \pi(X_i, Y_i, z(M, W^{i-1})) - o(1) \end{aligned} \quad (132)$$

$$\stackrel{(b)}{=} \min_{z(u)} \mathbb{E}_{\widehat{Q}} \pi(X_i, Y_i, z(u_i(M))) - o(1) \quad (133)$$

$$\stackrel{(c)}{=} \min_{z(u)} \mathbb{E} \pi(X, Y, z(U)) - o(1) \quad (134)$$

$$\geq \Pi - o(1). \quad (135)$$

Step (a) uses the first part of Lemma 6 along with Property 2b of total variation. Step (b) follows from Lemma 1 because under $\widehat{Q}_{MW^n X_i Y_i}$, the following markov chain holds:

$$(X_i, Y_i) - u_i(M) - (M, W^{i-1}). \quad (136)$$

Step (c) is due to the second part of Lemma 6 and Property 2b of total variation. This completes the analysis of the MIN payoff criterion.

E. Analysis of the WHP payoff criterion

Without loss of generality, we restrict attention to those distributions P_{UVXYW} that satisfy

$$P_{XY}(x, y) > 0 \implies \pi(x, y, z) > -\infty, \forall x, y, z. \quad (137)$$

Otherwise, $\min_z \mathbb{E} \pi(X, Y, z) = -\infty$ and the region in Theorem 1 is trivial.

The analysis will take place over sub-blocks of length $k = \lfloor \alpha n \rfloor$ rather than over the full block. For ease of presentation, we assume that $\lfloor \alpha n \rfloor = \alpha n$ and that k divides n evenly; the analysis is readily adjusted when this is not the case. We first fix some notation for handling sub-blocks. Denote the indices of the j th sub-block by the set $\mathcal{B}_{(j)}$:

$$\mathcal{B}_{(j)} = \{jk, jk + 1, \dots, (j+1)k - 1\}, \quad j \in [1/\alpha]. \quad (138)$$

Furthermore, denote the first t indices of sub-block $\mathcal{B}_{(j)}$ by $\mathcal{B}_{(j)}^t$; for example, $\mathcal{B}_{(j)}^1 = j$ and $\mathcal{B}_{(j)}^k = \mathcal{B}_{(j)}$. Some more notation: denote the adversary’s optimal reconstruction sequence by $\{Z_i^*\}_{i=1}^n$ and, for brevity, define

$$\rho \triangleq \min_{z(u)} \mathbb{E} \pi(X, Y, z(U)). \quad (139)$$

Let $\Pi < \rho$ and $\varepsilon = \rho - \Pi$. To prove achievability under the WHP criterion, we claim that it is enough to show that, for all $j \in [1/\alpha]$,

$$\lim_{k \rightarrow \infty} \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P}_{\widehat{Q}} \left[\frac{1}{k} \sum_{i \in \mathcal{B}_{(j)}} \pi(X_i, Y_i, Z_i^*) < \rho - \varepsilon \right] = 0, \quad (140)$$

where $\widehat{Q}_{MW^n X_{\mathcal{B}_{(j)}} Y_{\mathcal{B}_{(j)}}}$ is given in Lemma 5. Indeed, if this

is true, then we can write

$$\begin{aligned} & \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P} \left[\frac{1}{n} \sum_{i=1}^n \pi(X_i, Y_i, Z_i^*) \geq \Pi - \varepsilon \right] \\ & \geq \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P} \left[\frac{1}{n} \sum_{i=1}^n \pi(X_i, Y_i, Z_i^*) \geq \rho - \varepsilon \right] \end{aligned} \quad (141)$$

$$\geq \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P} \left[\bigcap_j \left\{ \frac{1}{k} \sum_{i \in \mathcal{B}^{(j)}} \pi(X_i, Y_i, Z_i^*) \geq \rho - \varepsilon \right\} \right] \quad (142)$$

$$= 1 - \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P} \left[\bigcup_j \left\{ \frac{1}{k} \sum_{i \in \mathcal{B}^{(j)}} \pi(X_i, Y_i, Z_i^*) < \rho - \varepsilon \right\} \right] \quad (143)$$

$$\stackrel{(a)}{\geq} 1 - \sum_{j=1}^{1/\alpha} \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P} \left[\frac{1}{k} \sum_{i \in \mathcal{B}^{(j)}} \pi(X_i, Y_i, Z_i^*) < \rho - \varepsilon \right] \quad (144)$$

$$\stackrel{(b)}{=} 1 - \sum_{j=1}^{1/\alpha} \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P}_{\bar{Q}} \left[\frac{1}{k} \sum_{i \in \mathcal{B}^{(j)}} \pi(X_i, Y_i, Z_i^*) < \rho - \varepsilon \right] - o(1) \quad (145)$$

$$\stackrel{(c)}{=} 1 - o(1). \quad (146)$$

Step (a) uses a union bound. Step (b) is due to Lemma 5 and the definition of total variation. Step (c) follows from the hypothesis in (140) and the fact that $1/\alpha$ is a constant that does not grow with n .

We now show that (140) holds for all $j \in [1/\alpha]$. Since our analysis is the same for all sub-blocks, we drop the subscript on $\mathcal{B}^{(j)}$ and simply consider an arbitrary sub-block \mathcal{B} of size k .

We cannot use the standard law of large numbers to show (140) because the dependence of Z_i^* on (M, W^{i-1}) implies that the random variables $\{\pi(X_i, Y_i, Z_i^*)\}_{i \in \mathcal{B}}$ are not mutually independent. Instead, we condition on $U^n(M)$ and use a martingale argument.

For simplicity, denote $U^n(M)$ by \bar{U}^n . Let $\{S_t\}_{t \in \mathcal{B}}$ be defined by

$$S_t \triangleq \sum_{i \in \mathcal{B}^t} (\pi(X_i, Y_i, Z_i^*) - \rho_i(\bar{U}^n)), \quad (147)$$

where

$$\rho_i(\bar{U}^n) \triangleq \min_{z(u^n)} \mathbb{E}_{\bar{Q}}[\pi(X_i, Y_i, z) | \bar{U}^n]. \quad (148)$$

We claim that, conditioned on \bar{U}^n , S_t is a submartingale, i.e.,

$$\mathbb{E}_{\bar{Q}}[S_t | S^{t-1}, \bar{U}^n] \geq S_{t-1}, \quad \forall t \in \mathcal{B}. \quad (149)$$

To verify the claim, first observe that the definition of S_t gives

$$\begin{aligned} \mathbb{E}_{\bar{Q}}[S_t | S^{t-1}, \bar{U}^n] &= S_{t-1} + \mathbb{E}_{\bar{Q}}[\pi(X_t, Y_t, Z_t^*) | S^{t-1}, \bar{U}^n] \\ &\quad - \rho_t(\bar{U}^n). \end{aligned} \quad (150)$$

Moreover, for each $t \in \mathcal{B}$, we have

$$\begin{aligned} & \mathbb{E}_{\bar{Q}}[\pi(X_t, Y_t, Z_t^*) | S^{t-1}, \bar{U}^n] \\ & \geq \min_{z(m, w^{t-1}, u^n, s^{t-1})} \mathbb{E}_{\bar{Q}}[\pi(X_t, Y_t, z(M, W^{t-1})) | S^{t-1}, \bar{U}^n] \end{aligned} \quad (151)$$

$$\stackrel{(a)}{=} \min_{z(u^n, s^{t-1})} \mathbb{E}_{\bar{Q}}[\pi(X_t, Y_t, z) | S^{t-1}, \bar{U}^n] \quad (152)$$

$$\stackrel{(b)}{=} \min_{z(u^n)} \mathbb{E}_{\bar{Q}}[\pi(X_t, Y_t, z) | \bar{U}^n] \quad (153)$$

$$= \rho_t(\bar{U}^n). \quad (154)$$

Step (a) follows by invoking Lemma 1 after noting that under \bar{Q} we have the markov chain

$$(X_t, Y_t) - (\bar{U}^n, S^{t-1}) - (M, W^{t-1}). \quad (155)$$

Step (b) follows from the markov chain

$$(X_t, Y_t) - \bar{U}^n - S^{t-1}. \quad (156)$$

Thus, conditioned on \bar{U}^n , we see that S_t is a submartingale. By Doob's decomposition theorem, we can write $S_t = M_t + A_t$, where M_t is a martingale (conditioned on \bar{U}^n) and A_t is an increasing process with $A_1 = 0$. Therefore, conditioning on \bar{U}^n , we have

$$\begin{aligned} & \mathbb{P}_{\bar{Q}} \left[\frac{1}{k} \sum_{i \in \mathcal{B}} \pi(X_i, Y_i, Z_i^*) < \frac{1}{k} \sum_{i \in \mathcal{B}} \rho_i(\bar{U}^n) - \varepsilon \mid \bar{U}^n \right] \\ & = \mathbb{P}_{\bar{Q}}[S_k < -k\varepsilon | \bar{U}^n] \end{aligned} \quad (157)$$

$$\leq \mathbb{P}_{\bar{Q}}[M_k < -k\varepsilon | \bar{U}^n] \quad (158)$$

$$= \mathbb{P}_{\bar{Q}}[M_k - \mathbb{E}_{\bar{Q}}[M_k] < -k\varepsilon - \mathbb{E}_{\bar{Q}}[M_k] | \bar{U}^n] \quad (159)$$

$$\stackrel{(a)}{\leq} \frac{\text{Var}_{\bar{Q}}(M_k | \bar{U}^n)}{(k\varepsilon + \mathbb{E}_{\bar{Q}}[S_1])^2}, \quad (160)$$

where (a) follows from Chebyshev's inequality. Now we recursively bound the variance of M_k (conditioned on \bar{U}^n) by writing

$$\begin{aligned} & \text{Var}(M_k | \bar{U}^n) \\ & \stackrel{(a)}{=} \text{Var}(\mathbb{E}[M_k | M^{k-1}, \bar{U}^n]) \end{aligned} \quad (161)$$

$$+ \mathbb{E}[\text{Var}(M_k | M^{k-1}, \bar{U}^n)] \quad (162)$$

$$\leq \text{Var}(\mathbb{E}[M_k | M^{k-1}, \bar{U}^n]) + O(1) \quad (163)$$

$$= \text{Var}(M_{k-1} | \bar{U}^n) + O(1). \quad (164)$$

Step (a) uses the law of total variance. The recursion implies $\text{Var}_{\bar{Q}}(M_k | \bar{U}^n) \in O(k)$, which, together with (160), shows

$$\lim_{k \rightarrow \infty} \mathbb{P}_{\bar{Q}} \left[\frac{1}{k} \sum_{i \in \mathcal{B}} \pi(X_i, Y_i, Z_i^*) < \frac{1}{k} \sum_{i \in \mathcal{B}} \rho_i(\bar{U}^n) - \varepsilon \mid \bar{U}^n \right] = 0. \quad (165)$$

Since this convergence is uniform for all \bar{U}^n , we can take the expectation over random codebooks to get

$$\lim_{k \rightarrow \infty} \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P}_{\bar{Q}} \left[\frac{1}{k} \sum_{i \in \mathcal{B}} \pi(X_i, Y_i, Z_i^*) < \frac{1}{k} \sum_{i \in \mathcal{B}} \rho_i(\bar{U}^n) - \varepsilon \right] = 0. \quad (166)$$

Continuing, notice that $\rho_i(\bar{U}^n)$ can be written as

$$\rho_i(\bar{U}^n) = \min_z \mathbb{E}_{\bar{Q}}[\pi(X_i, Y_i, z) | \bar{U}^n] \quad (167)$$

$$= \min_z \mathbb{E}_{\bar{Q}}[\pi(X_i, Y_i, z) | \bar{U}_i] \quad (168)$$

$$\triangleq \rho(\bar{U}_i) \quad (169)$$

because of the markov chain $(X_i, Y_i) - \bar{U}_i - \bar{U}^n$ that holds under \hat{Q} . Furthermore, the expected value of $\rho(\bar{U}_i)$ is

$$\mathbb{E}_{\mathcal{C}^{(n)}} \rho(\bar{U}_i) = \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{E}_{\hat{Q}}[\pi(X_i, Y_i, z) | \bar{U}_i] \quad (170)$$

$$= \min_{z(u)} \mathbb{E}_{\mathcal{C}^{(n)}} \pi(X_i, Y_i, z(\bar{U}_i)) \quad (171)$$

$$\stackrel{(a)}{=} \min_{z(u)} \mathbb{E}_{\mathcal{C}^{(n)}} \pi(X, Y, z(U)) \quad (172)$$

$$= \rho \quad (173)$$

where step (a) is due to the fact (readily verified) that $\mathbb{E}_{\mathcal{C}^{(n)}} Q_{X_i Y_i U_i(M)} = P_{XYU}$. Therefore, because \bar{U}^n is i.i.d. according to P_U (in expectation over the random codebooks), we can invoke the law of large numbers to get

$$\lim_{k \rightarrow \infty} \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P}_{\hat{Q}} \left[\frac{1}{k} \sum_{i \in \mathcal{B}} \rho(\bar{U}_i) > \rho - \varepsilon \right] = 1. \quad (174)$$

This, together with (166), yields

$$\lim_{k \rightarrow \infty} \mathbb{E}_{\mathcal{C}^{(n)}} \mathbb{P}_{\hat{Q}} \left[\frac{1}{k} \sum_{i \in \mathcal{B}} \pi(X_i, Y_i, Z_i^*) < \rho - 2\varepsilon \right] = 0, \quad (175)$$

completing the proof of (140). Finally, we invoke Shannon's random coding argument to ensure the existence of a codebook that satisfies the payoff criterion. This concludes the achievability proof of the WHP payoff criterion.

F. Discussion: Optimal encoding produces artificial noise

The optimal encoding and decoding scheme designed in this section produces an effect that is worth investigating outside of this particular context of rate-distortion theory for secrecy systems. In particular, consider the most pessimistic disclosure assumption, that $W = (X, Y)$. In this case, the communication system effectively corrupts the i.i.d. information signal X^n with noise by synthesizing a memoryless broadcast channel, with the information source X^n as input, actions at the intended receiver Y^n as one output, and a sequence U^n as the other output observed by the adversary. The synthesis is accurate in a particular sense relevant to secrecy. That is, the communication system, which uses public message M and secret key K to facilitate coordination, synthesizes memoryless noise characterized by $P_{YU|X}$ by producing a distribution on (X^n, Y^n, M) such that $P_{X^n Y^n | M}$ closely approximates $\prod_{i=1}^n P_{X_i Y_i | U_i(M)}$ for a set of statistically typical $u^n(M)$ sequences. This behavior is revealed by $\hat{Q}_{M X^n Y^n}$ in (121), which the proof shows to converge to the induced joint distribution of the system in the limit of large n .

Let us now consider why this might be an operationally meaningful criterion for synthesizing noise in a secrecy setting. Consider an adversary who actually does observe a noise-corrupted version of the information signal, such as one of the outputs of a broadcast channel. As in any probabilistic situation, rational behavior is based on the posterior distribution of the state of the universe given what is known to the individual. In this situation that means $P_{X^n, Y^n | U^n}$ will dictate the adversary's optimal behavior, regardless of the objective that the adversary is trying to accomplish. Therefore, a communication system that mimics $P_{X^n, Y^n | U^n}$ will elicit the same behavior by an adversary for the same observed U^n sequence

as would occur if the noisy channel was genuine. Furthermore, if the observed U^n sequence is statistically representative⁹ of true noisy observations, then the communication system performance in the presence of an adversary will be equivalent to the memoryless broadcast channel that it mimics.

For comparison, consider the work of Winter in [20]. Although the communication setting and results in [20] are quite different from ours in that the setting does not have an information source provided by nature, our proof and methods for achievability bear resemblance. There, he considers a distribution on a triple of variables (X, Y, U) and a communication system that generates correlated random variables X^n and Y^n at two different nodes using communication and secret key in the presence of an adversary. For the sake of comparison, imagine Y^n as a noisy version of X^n . The secrecy criterion in that work is very strong, requiring that the public message reveal no more about the sequences X^n and Y^n than the correlated sequence U^n would, in the sense that M is stochastically degraded from U^n with respect to (X^n, Y^n) . This is stronger than the secrecy criterion we gave in the previous paragraphs, requiring more communication resources as a consequence. However, the noise synthesis achieved by the communication system of this section, even with the weaker secrecy performance implied by (121), has the same compelling operational significance—an adversary can gain no more advantage from the eavesdropped message than they could by observing the correlated U^n sequence.

VII. CONVERSE PROOF

It is enough to prove the converse to Theorem 1 for just the AVG payoff criterion, since it is the weakest of the criteria. We further weaken the conditions by allowing Node B causal access to Nodes A and C (i.e., we permit decoders of the form $\{P_{Y_i | M K X^{i-1} Z^{i-1}}\}_{i=1}^n$). We will see that this allowance does not increase the payoff.

Fix a source distribution P_X , a payoff function $\pi(x, y, z)$, and causal disclosure channels $P_{W_x | X}$ and $P_{W_y | Y}$. For ease of presentation, denote the pair (W_x^n, W_y^n) by W^n . Next, let J be an auxiliary random variable drawn uniformly from $[n]$, independently of (X^n, Y^n, W^n, M, K) . Define the following random variables:

$$X = X_J \quad (176)$$

$$Y = Y_J \quad (177)$$

$$Z = Z_J \quad (178)$$

$$(W_x, W_y) = W_J \quad (179)$$

$$U = (M, W^{J-1}, J) \quad (180)$$

$$V = K. \quad (181)$$

With these choices, it can be verified that

$$W_x - X - (U, V) - Y - W_y \quad (182)$$

$$X \sim P_X \quad (183)$$

$$W_x | X \sim P_{W_x | X} \quad (184)$$

$$W_y | Y \sim P_{W_y | Y} \quad (185)$$

⁹Exact characterization of this depends on the specific objectives of the communication system.

The following properties of $P_{MKX^nY^nW^n}$ can also be verified:

$$X^n \perp K \quad (186)$$

$$X_J - (M, K, X^{J-1}, J) - W^{J-1} \quad (187)$$

$$X_J \perp J. \quad (188)$$

Let (R, R_0, Π) be an achievable triple. We first have

$$nR \geq H(M) \quad (189)$$

$$\geq H(M|K) \quad (190)$$

$$\geq I(X^n; M|K) \quad (191)$$

$$\stackrel{(a)}{=} I(X^n; M, K) \quad (192)$$

$$= \sum_{j=1}^n I(X_j; M, K|X^{j-1}) \quad (193)$$

$$= \sum_{j=1}^n I(X_j; M, K, X^{j-1}) \quad (194)$$

$$\stackrel{(b)}{=} \sum_{j=1}^n I(X_j; M, K, X^{j-1}, W^{j-1}) \quad (195)$$

$$\geq \sum_{i=1}^n I(X_i; M, K, W^{i-1}) \quad (196)$$

$$\stackrel{(c)}{=} nI(X_J; M, K, W^{J-1}, J) \quad (197)$$

$$= nI(X; U, V), \quad (198)$$

where (a), (b), and (c) follow from (186), (187), and (188). Next, we have

$$nR_0 \geq H(K) \quad (199)$$

$$\geq H(K|M) \quad (200)$$

$$\geq I(W^n; K|M) \quad (201)$$

$$\geq \sum_{j=1}^n I(W_j; K|M, W^{j-1}) \quad (202)$$

$$\stackrel{(a)}{=} nI(W_J; K|M, W^{J-1}, J) \quad (203)$$

$$= nI(W; V|U), \quad (204)$$

where (a) follows from (188). Finally, we have

$$\Pi \leq \min_{z(m, w^{j-1}, j)} \mathbb{E} \frac{1}{n} \sum_{j=1}^n \pi(X_j, Y_j, z(M, W^{j-1}, j)) \quad (205)$$

$$= \min_{z(m, w^{j-1}, j)} \mathbb{E} \left[\mathbb{E}[\pi(X_J, Y_J, z(M, W^{J-1}, J)) | J] \right] \quad (206)$$

$$= \min_{z(m, w^{j-1}, j)} \mathbb{E} \pi(X_J, Y_J, z(M, W^{J-1}, J)) \quad (207)$$

$$= \min_{z(u)} \mathbb{E} \pi(X, Y, z(U)). \quad (208)$$

It remains to bound the cardinality of \mathcal{U} and \mathcal{V} , which is straightforward from the standard support lemma (e.g., [21]). Note that the set of markov distributions forms a compact, connected set. To bound \mathcal{U} , it suffices to have $|\mathcal{X}|-1$ elements to preserve P_X and 3 more elements to preserve $H(X|U, V)$, $I(W; V|U)$, and $\min_{z(u)} \mathbb{E} \pi(X, Y, z(U))$. To bound \mathcal{V} , it suffices to have $|\mathcal{X}||\mathcal{Y}||\mathcal{U}|-1$ elements to preserve P_{XYU} and 2 more elements to preserve $H(X|U, V)$ and $H(W|U, V)$.

VIII. OTHER FORMS OF DISCLOSURE

In this section, we consider several relevant scenarios that are not directly subsumed by Theorem 1, but that can be solved by modifying the proof slightly. Throughout, we denote (W_x^n, W_y^n) by W^n . Whereas previously we assumed that the eavesdropper has access to causal disclosure W^{i-1} , now we consider three other types of disclosure: W_i , W^i , and W^n . It turns out that the regions corresponding to W^i and W^n are the same.

Theorem 3: Fix P_X , $\pi(x, y, z)$ and disclosure channels $P_{W_x|X}$ and $P_{W_y|Y}$. If W_i is disclosed instead of W^{i-1} , then the rate-payoff region for all three payoff criteria is equal to

$$\bigcup_{P_{Y|X}} \left\{ (R, R_0, \Pi) : \begin{array}{l} R \geq I(X; Y) \\ R_0 \geq 0 \\ \Pi \leq \min_{z(w_x, w_y)} \mathbb{E} \pi(X, Y, z(W_x, W_y)) \end{array} \right\}. \quad (209)$$

Proof: The proof of achievability is very similar to that of Section VI. Define the random codebook, encoder, decoder, and $Q_{X^n M K Y^n W^n}$ in the same way, but set $U = \emptyset$ and $V = Y$ throughout. Lemma 4 ensures that the system-induced distribution is approximated by Q since $R > I(X; Y)$. Instead of the property in Lemma 5, the desired property of Q is now

$$Q_{M X_B Y_B W_B} \approx Q_M \cdot \left(\prod_{i \in \mathcal{B}} P_{X_i Y_i W_i} \right). \quad (210)$$

The soft covering lemma can be invoked to show that this property holds if the rate of secret key satisfies

$$R_0 > \limsup_{n \rightarrow \infty} \frac{1}{n} i_Q(X_B Y_B W_B; Y^n) = 0. \quad (211)$$

Thus, under Q , the message M is approximately independent of (X_i, Y_i, W_i) and the eavesdropper's best estimate of (X_i, Y_i) only depends on his observation of the disclosure W_i . The payoff analysis of Section VI is straightforward to modify accordingly.

To prove the converse, it is first straightforward to bound R and R_0 . To bound Π , define $(W_x, W_y) = W_J$, where $J \sim \text{Unif}(n)$, and write

$$\Pi \leq \min_{z(m, w_j, j)} \mathbb{E} \frac{1}{n} \sum_{j=1}^n \pi(X_j, Y_j, z(M, W_j, j)) \quad (212)$$

$$\leq \min_{z(w_j)} \mathbb{E} \frac{1}{n} \sum_{j=1}^n \pi(X_j, Y_j, z(W_j)) \quad (213)$$

$$= \min_{z(w)} \mathbb{E} \pi(X_J, Y_J, z(W_J)) \quad (214)$$

$$= \min_{z(w_x, w_y)} \mathbb{E} \pi(X, Y, z(W_x, W_y)). \quad (215)$$

Theorem 4: Fix P_X , $\pi(x, y, z)$ and disclosure channels $P_{W_x|X}$ and $P_{W_y|Y}$. If W^n or W^i is disclosed instead of W^{i-1} , then the rate-payoff region for all three payoff criteria is equal

to

$$\bigcup \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq I(X; U, V) \\ R_0 \geq I(W_x, W_y; V|U) \\ \Pi \leq \min_{z(u, w_x, w_y)} \mathbb{E} \pi(X, z(U, W_x, W_y)) \end{array} \right\}, \quad (216)$$

where the union is taken over all markov chains

$$W_x - X - (U, V) - Y - W_y. \quad (217)$$

Proof: For the proof of achievability, suppose W^n is disclosed. The proof is almost exactly the same as in Section VI. The code and the rates are identical, as is the definition of the approximating distribution Q . Notice that under \widehat{Q} (defined in Lemma 5), the following markov chain holds for all $i \in [n]$:

$$(X_i, Y_i) - (U_i(M), W_i) - (M, W^n). \quad (218)$$

Thus, the eavesdropper's best strategy only depends on $(U_i(M), W_i)$; the rest of the disclosure of W^n is rendered useless. To adjust the analysis of the payoff criteria, simply use markov relations similar to the one in (218).

To show the converse proof, suppose that only W^i is disclosed. The proof follows arguments similar to those in Section VII, with exactly the same identification of random variables. ■

IX. CAUSAL DISCLOSURE WITH DELAY

In this section, we consider the effects of assuming that the adversary has delayed causal access to the system behavior. In other words, we replace causal disclosure W^{i-1} with W^{i-d} , $d > 1$. Surprisingly, this has a major effect on relaxing the amount of secret key required to maintain secrecy. We establish an inner and outer bound on the corresponding rate-payoff region and give an example in which the bounds meet.¹⁰ Using the bounds, we further show that if lossless communication is required, the minimum rate of secret key needed to ensure a given level of payoff is on the order of $1/d$.

A. Inner and outer bound

Theorem 5 (Inner bound, causal disclosure with delay d): Fix P_X , $\pi(x, y, z)$, and causal disclosure channels $P_{W_x|X}$ and $P_{W_y|Y}$. Let \mathcal{R}_d denote the closure of achievable (R, R_0, Π) when the causal disclosure has delay $d \geq 1$. Then

$$\mathcal{R}_d \supseteq \bigcup \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq \frac{1}{d} I(X^d; U, V) \\ R_0 \geq \frac{1}{d} I(W_x^d W_y^d; V|U) \\ \Pi \leq \min_{z(u)} \mathbb{E} \left[\frac{1}{d} \sum_{j=1}^d \pi(X_j, Y_j, z(U)) \right] \end{array} \right\}, \quad (219)$$

where the union is taken over all markov chains

$$W_x^d - X^d - (U, V) - Y^d - W_y^d \quad (220)$$

¹⁰Numerical investigation reveals that the bounds are not tight in general.

in which

$$P_{X^d W_x^d} = \prod_{j=1}^d P_X P_{W_x|X} \quad (221)$$

and

$$P_{W_y^d|Y^d} = \prod_{j=1}^d P_{W_y|Y}. \quad (222)$$

Proof: For simplicity, we present the proof for $d = 2$. Denote (W_x^n, W_y^n) by W^n . The idea is to transform the problem into one involving delay $d = 1$ so that we can apply Theorem 1. To that end, we first treat the source X^n as an i.i.d. sequence $\widetilde{X}^{\frac{n}{2}}$ of super-symbols of length 2 by defining

$$\widetilde{X}_i = (X_{2i-1}, X_{2i}), \quad i = 1, 2, \dots, n/2. \quad (223)$$

Similarly, treat Y^n and W^n as sequences of super-symbols by appropriately defining $\widetilde{Y}^{\frac{n}{2}}$ and $\widetilde{W}^{\frac{n}{2}}$. Under this definition, observe that at steps $i = 2, 4, \dots, n$ the adversary has access to $\widetilde{W}^{i-1} = W^{i-2}$. Suppose that at steps $i = 1, 3, \dots, n$ we disclose additional information W_{i-1} to the adversary. Now the causal disclosure to the adversary is exactly \widetilde{W}^{i-1} for all $i \in [n]$. Note that supplying extra information to the adversary can only reduce the achievable region.

To complete the transformation, define a payoff function $\widetilde{\pi} : \mathcal{X}^2 \times \mathcal{Y}^2 \times \mathcal{Z}^2 \rightarrow \mathbb{R}$ by

$$\widetilde{\pi}(x^2, y^2, z^2) = \sum_{j=1}^2 \pi(x_j, y_j, z_j). \quad (224)$$

If $(\widetilde{R}, \widetilde{R}_0, \widetilde{\Pi})$ is an achievable triple for this transformed problem, then $(\widetilde{R}/2, \widetilde{R}_0/2, \widetilde{\Pi}/2)$ is an achievable triple for the delayed causal disclosure problem with $d = 2$. By applying Theorem 1, we obtain the region in (219) for $d = 2$. ■

Theorem 6 (Outer bound, causal disclosure with delay d): Fix P_X , $\pi(x, y, z)$, and causal disclosure channels $P_{W_x|X}$ and $P_{W_y|Y}$. Let \mathcal{R}_d denote the closure of achievable (R, R_0, Π) when the causal disclosure has delay $d \geq 1$. Then

$$\mathcal{R}_d \subseteq \bigcup \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq I(X; U, V) \\ R_0 \geq \frac{1}{d} I(W_x W_y; V|U) \\ \Pi \leq \min_{z(u)} \mathbb{E} \pi(X, Y, z(U)) \end{array} \right\}, \quad (225)$$

where the union is taken over all markov chains

$$W_x - X - (U, V) - Y - W_y. \quad (226)$$

Proof: The key to the proof is the following lemma.

Lemma 7: For arbitrary random variables (X^n, Y) , it holds that

$$d \cdot I(X^n; Y) \geq \sum_{i=1}^n I(X_i; Y|X^{i-d}). \quad (227)$$

Proof of Lemma 7:

$$\begin{aligned} d \cdot I(X^n; Y) &= \sum_{j=1}^d I(X^n; Y) \\ &\geq \sum_{j=1}^d I(X_j^{n - ((n-j) \bmod d)}; Y) \end{aligned} \quad (228)$$

$$\stackrel{(a)}{=} \sum_{j=1}^d \sum_{i \in [n], i \geq j, i \equiv j \pmod{d}} I(X_{i-d+1}^i; Y | X^{i-d}) \quad (230)$$

$$= \sum_{i=1}^n I(X_{i-d+1}^i; Y | X^{i-d}) \quad (231)$$

$$\geq \sum_{i=1}^n I(X_i; Y | X^{i-d}). \quad (232)$$

Step (a) uses the chain rule for mutual information on each of the d terms. ■

The converse steps of Section VII can now be modified by defining $U = (M, W^{J-d}, J)$.

First, bound R by writing

$$nR \geq H(M) \quad (233)$$

⋮

$$\geq nI(X_J; M, K, W^{J-1}, J) \quad (234)$$

$$\geq nI(X_J; M, K, W^{J-d}, J) \quad (235)$$

$$= nI(X; U, V). \quad (236)$$

Next, bound R_0 by writing

$$d \cdot nR \geq d \cdot H(M) \quad (237)$$

⋮

$$\geq d \cdot I(W^n; K | M) \quad (238)$$

$$\stackrel{(a)}{\geq} \sum_{j=1}^n I(W_j; K | M, W^{j-d}) \quad (239)$$

$$= nI(W_J; K | M, W^{J-d}, J) \quad (240)$$

$$= nI(W; V | U), \quad (241)$$

where (a) uses Lemma 7. Finally, Π can be bounded in the manner of Section VII. ■

B. Lossless communication

We now specialize the inner and outer bound to the setting in which lossless communication is required and X^{i-d} is disclosed. In this regime, we are able to show explicitly how delay affects the tradeoff between rate of secret key and payoff.

Theorem 7: Fix P_X and $\pi(x, z)$. Let \mathcal{R}_d denote the closure of achievable (R, R_0, Π) for the case of lossless communication and causal disclosure X^{i-d} , $d \geq 1$. Let $R_d(\Pi)$ denote

the key-payoff boundary of \mathcal{R}_d . First, we have

$$\mathcal{R}_d \supseteq \bigcup_{X^d \sim \prod_{j=1}^d P_X} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq H(X) \\ R_0 \geq \frac{1}{d} H(X^d | U) \\ \Pi \leq \min_{z(u)} \mathbb{E} \left[\frac{1}{d} \sum_{j=1}^d \pi(X_j, z(U)) \right] \end{array} \right\} \quad (242)$$

and

$$\mathcal{R}_d \subseteq \bigcup_{\substack{P_{XU}: \\ X \sim P_X}} \left\{ \begin{array}{l} (R, R_0, \Pi) : \\ R \geq H(X) \\ R_0 \geq \frac{1}{d} H(X | U) \\ \Pi \leq \min_{z(u)} \mathbb{E} \pi(X, z(U)) \end{array} \right\}. \quad (243)$$

Furthermore, for all Π ,

$$R_d(\Pi) = \Theta\left(\frac{1}{d}\right). \quad (244)$$

Proof: To establish the inner bound on \mathcal{R}_d , first recall the characterization of \mathcal{R}_1 given in Corollary 3. Using the same arguments as the proof of Theorem 5, we can transform the problem with delay $d > 1$ into one involving delay $d = 1$ and invoke Corollary 3 on the new problem. Upon noting that $\frac{1}{d} H(X^d) = H(X)$ when $X^d \sim \prod P_X$, this technique gives the achievable region in (242).

To establish the outer bound, let (R, R_0, Π) be an achievable triple. The bound $R \geq H(X)$ is due to the lossless source coding theorem. To bound R_0 and Π , let J be uniformly distributed on $[n]$ and define $U = (M, X^{J-d}, J)$ and $X = X_J$. Then, we have

$$nR_0 \geq nH(K) \quad (245)$$

$$\geq nI(X^n; K | M) \quad (246)$$

$$= nH(X^n | M) - nH(X^n | K, M) \quad (247)$$

$$\stackrel{(a)}{=} nH(X^n | M) - n \cdot o(1) \quad (248)$$

$$\stackrel{(b)}{\geq} n \cdot \frac{1}{d} \sum_{j=1}^n H(X_j | X^{j-d}, M) - n \cdot o(1) \quad (249)$$

$$= n \cdot \frac{1}{d} H(X_J | X^{J-1}, M, J) - n \cdot o(1) \quad (250)$$

$$= n \cdot \frac{1}{d} H(X | U) - n \cdot o(1). \quad (251)$$

Step (a) uses Fano's inequality, and step (b) follows from Lemma 7 (by setting $Y = X^n$ and conditioning on M). It is straightforward to bound Π in the manner of Section VII.

From the outer bound in (243), we see that $R_d(\Pi) \geq \frac{1}{d} R_1(\Pi)$. It remains to show that $R_d(\Pi) \leq c \cdot \frac{1}{d}$ for some constant c ; we do this via (242). First, let $X^d \sim \prod_{j=1}^d P_X$. Let $K \sim \text{Unif}(\mathcal{X})$ be independent of X^d and define

$$U \triangleq (X_1 \oplus K, X_2 \oplus K, \dots, X_d \oplus K), \quad (252)$$

where \oplus indicates addition modulo \mathcal{X} . With this choice of U , we have

$$H(X_i | X_j, U) = 0, \forall i, j \in [d] \quad (253)$$

and

$$X_j \perp U, \forall j \in [d]. \quad (254)$$

Therefore, we can write

$$\frac{1}{d}H(X^d|U) = \frac{1}{d}\sum_{i=1}^d H(X_i|X^{i-1}, U) \quad (255)$$

$$\stackrel{(a)}{=} \frac{1}{d}H(X_1|U) \quad (256)$$

$$\stackrel{(b)}{=} \frac{1}{d}H(X), \quad (257)$$

where (a) and (b) follow from (253) and (254), respectively. Moreover, we have

$$\min_{z(u)} \mathbb{E} \left[\frac{1}{d} \sum_{j=1}^d \pi(X_j, z(U)) \right] \quad (258)$$

$$= \frac{1}{d} \sum_{j=1}^d \min_{z(u)} \mathbb{E} \pi(X_j, z(U)) \quad (259)$$

$$\stackrel{(a)}{=} \min_z \mathbb{E} \pi(X, z) \quad (260)$$

$$\triangleq \pi_{\max}, \quad (261)$$

where (a) follows from Lemma 1 and (254).

By selecting U according to (252), we have shown that the inner bound in (242) contains the point $(R_0, \Pi) = (\frac{1}{d}H(X), \pi_{\max})$; therefore, $(\frac{1}{d}H(X), \pi_{\max}) \in \mathcal{R}_d$. Since π_{\max} is the maximum possible payoff, this implies $R_d(\Pi) \leq \frac{1}{d}H(X)$, completing the proof of (244). ■

C. Example in which the bounds meet

In the preceding proof, we demonstrated that the point $(R_0, \Pi) = (\frac{1}{d}H(X), \pi_{\max})$ is in the region (242) and is therefore achievable. If we choose the source distribution to be $P_X \sim \text{Bern}(1/2)$, then from Theorem 2 (which gives us $R_1(\Pi)$) and the convexity of the rate-payoff region, it is clear that $R_d(\Pi) \leq \frac{1}{d}R_1(\Pi)$. Conversely, the outer bound in (243) directly gives $R_d \geq \frac{1}{d}R_1(\Pi)$.

X. CONCLUSION

This work has established a theory of secure source coding which characterizes the optimal use of communication and secret key to allow good reconstruction of the source by the intended receiver (who has access to the key) and force a poor reconstruction on any eavesdropper (without the secret key). The central contribution, presented in Theorem 1, gives a general information theoretic characterization of the achievable performance. The expression in the theorem makes use of two auxiliary variables which can be interpreted as information that is kept secure and information that is released publicly. In the case of lossless compression in Corollary 3, the optimal communication system can explicitly follow these implied steps, constructing two separate messages and focusing all of the security resource (i.e., the key) on only one.

An important component of the main result is the causal disclosure assumption depicted in Fig. 2, which was absent from Yamamoto's formulation of the problem in [11] and [12]. The causal disclosure empowers the eavesdropper with additional information and forces the communication system to resort to a more robust design for secure encoding, which

results in an innovative encoding and decoding scheme that sterilizes the causal disclosure.

The theorems in this work allow for an arbitrary but known disclosure channel to the eavesdropper. However, one could always take the most pessimistic approach and assume that the source X and the reconstruction Y are both fully disclosed (causally) to the eavesdropper. This leads to the strongest definition of secrecy in our model, and the optimal communication system for this setting has a simple and natural interpretation as producing synthetic noise, discussed in Section VI-F.

This work also identifies the rate-distortion tradeoffs without the causal disclosure assumption. The case of no disclosure (as in Yamamoto's model) is a special case of the main result and is addressed in Section III, along with a discussion of its fragility. Non-causal disclosure is the topic of Section VIII, which turns out to only be as empowering to the eavesdropper as causal disclosure.

The causal disclosure framework boasts some important unique properties aside from its operational interpretation as real-time reconstruction by the eavesdropper. In Section V we show that the traditional approach of measuring secrecy by normalized equivocation (rather than distortion) is in fact a special case of this framework by applying a particular log-loss distortion function. This connection only exists because of the causal disclosure assumption. Another property that arises is the need for a stochastic decoder, which suggests a duality with Wyner's wiretap channel [7] where a stochastic encoder is needed. Furthermore, this framework induces a rich tradeoff between the rate of secret key used and the distortion the system imposes upon an eavesdropper, while such a tradeoff does not occur in the absence of causal disclosure. These features suggest that causal disclosure is an appropriate base assumption for understanding rate-distortion theory for secrecy systems.

APPENDIX A

PROOF OF THEOREM 2

A. Supporting lemma

For each $x \in \mathcal{X}$, define $\mathcal{F}_n(x) \subseteq \Delta_{\mathcal{X}}$ by

$$\mathcal{F}_n(x) \triangleq \left\{ p \in \Delta_{\mathcal{X}} : p = \text{Unif}(\mathcal{A}) \text{ for some } \mathcal{A} \subseteq \mathcal{X}, \right. \\ \left. |\mathcal{A}| = n, \text{ and } p(x) = \max_{x'} p(x') \right\}, \quad (262)$$

and define $\mathcal{A}_n(x) \subseteq \Delta_{\mathcal{X}}$ by

$$\mathcal{A}_n(x) \triangleq \left\{ p \in \Delta_{\mathcal{X}} : p(x) = \max_{x'} p(x') \right. \\ \left. \text{and } p(x) \in \left[\frac{1}{n+1}, \frac{1}{n} \right] \right\}. \quad (263)$$

In words, $\mathcal{F}_n(x)$ is the set of probability mass functions on \mathcal{X} that are uniformly distributed on a subset of size n and whose largest mass occurs at x . Fig. 6 illustrates the definitions of $\mathcal{F}_n(x)$ and $\mathcal{A}_n(x)$ when $\mathcal{X} = \{1, 2, 3\}$.

The key to the proof of Theorem 2 is the following technical lemma.

Lemma 8: For a random variable X with distribution P_X , let \bar{x} and N be such that $P_X \in \mathcal{A}_N(\bar{x})$.

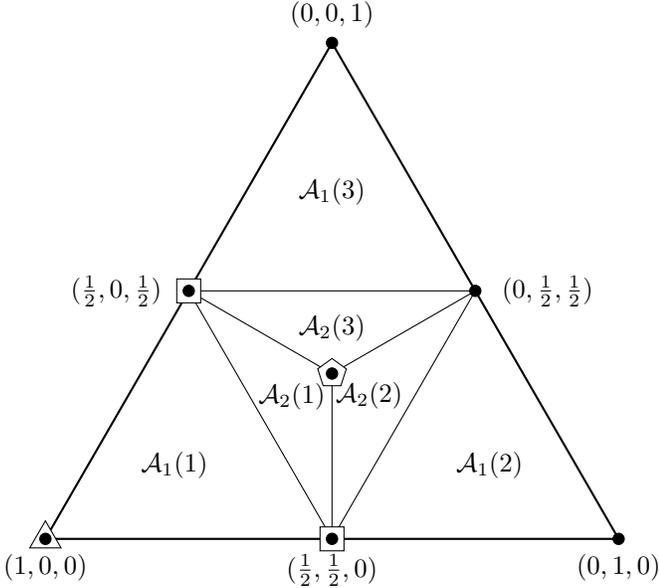


Fig. 6: The probability simplex $\Delta_{\mathcal{X}}$ for $\mathcal{X} = \{1, 2, 3\}$. The centroid is the distribution $(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$. Note that $\mathcal{F}_1(1) = \{\triangle\}$, $\mathcal{F}_2(1) = \{\square\}$, and $\mathcal{F}_3(1) = \{\diamond\}$.

- 1) There exists a random variable V , correlated with X , such that for all $v \in \mathcal{V}$,

$$P_{X|V=v} \in \mathcal{F}_N(\bar{x}) \cup \mathcal{F}_{N+1}(\bar{x}). \quad (264)$$

In other words, P_X can be written as a convex combination of distributions in $\mathcal{F}_N(\bar{x}) \cup \mathcal{F}_{N+1}(\bar{x})$.

- 2) Let $n \in [N]$. There exists a random variable V such that for all $v \in \mathcal{V}$,

$$P_{X|V=v} \in \bigcup_{x \in \mathcal{X}} \mathcal{F}_n(x). \quad (265)$$

In other words, for any $n \in [N]$, P_X can be written as a convex combination of distributions in $\bigcup_x \mathcal{F}_n(x)$.

Proof: Fix $\bar{x} \in \mathcal{X}$ and $n \in \mathbb{N}$, and define

$$\mathcal{F} \triangleq \mathcal{F}_n(\bar{x}) \cup \mathcal{F}_{n+1}(\bar{x}). \quad (266)$$

First, one can verify that $\mathcal{A}_n(\bar{x})$ is a convex set. Furthermore, it is well-known that every compact convex set is the convex hull of its extreme points. Thus, to prove part 1, it is enough to show that the set of extreme points of $\mathcal{A}_n(\bar{x})$ is equal to \mathcal{F} . Then any $p \in \mathcal{A}_n(\bar{x})$ can be written as a convex combination of the elements of \mathcal{F} .

The set of extreme points of a convex set \mathcal{C} is defined by

$$\begin{aligned} \text{extr}(\mathcal{C}) \triangleq \{p \in \mathcal{C} : \text{if } p = \theta q + (1 - \theta)r, q, r \in \mathcal{C}, \\ \theta \in (0, 1) \text{ then } p = q = r\}. \end{aligned} \quad (267)$$

We first show that $\mathcal{F} \subseteq \text{extr}(\mathcal{A}_n(\bar{x}))$. Let $p \in \mathcal{F}$, and let $q, r \in \mathcal{A}_n(\bar{x})$, $\theta \in (0, 1)$ be such that $q \neq p$, $r \neq p$, and

$$p = \theta q + (1 - \theta)r \quad (268)$$

If $p \in \mathcal{F}_n(\bar{x})$, then $p = q = r$ is clear because $q(x) \in [0, \frac{1}{n}]$ and $r(x) \in [0, \frac{1}{n}]$ for all $x \in \mathcal{X}$. On the other hand, suppose $p \in \mathcal{F}_{n+1}(\bar{x})$. Because $q, r \in \mathcal{A}_n(\bar{x})$ and $p(\bar{x}) = \frac{1}{n+1}$, we

have $q(\bar{x}) = r(\bar{x}) = \frac{1}{n+1}$. Thus, $q(x) \in [0, \frac{1}{n+1}]$ and $r(x) \in [0, \frac{1}{n+1}]$ for all $x \in \mathcal{X}$, and again $p = q = r$.

To show $\text{extr}(\mathcal{A}_n(\bar{x})) \subseteq \mathcal{F}$, we proceed by way of contradiction and suppose that $p \in \text{extr}(\mathcal{A}_n(\bar{x}))$ and $p \notin \mathcal{F}$. From $p \notin \mathcal{F}$, it holds that $p(x') \in (0, \frac{1}{n+1}) \cup (\frac{1}{n+1}, \frac{1}{n})$ for some $x' \in \mathcal{X}$. There are now three separate cases to consider depending on whether $p(\bar{x}) = \frac{1}{n+1}$, $p(\bar{x}) \in (\frac{1}{n+1}, \frac{1}{n})$, or $p(\bar{x}) = \frac{1}{n}$. For ease of exposition, we only consider $p(\bar{x}) = \frac{1}{n+1}$; the other two cases use a similar argument. Since $p(x') \leq p(\bar{x})$, we have $p(x') \in (0, \frac{1}{n+1})$. It follows that there must exist $x'' \neq x'$ such that $p(x'') \in (0, \frac{1}{n+1})$; otherwise, we would have

$$\sum_{x \in \mathcal{X}} p(x) = \frac{n}{n+1} + p(x') < 1 \quad (269)$$

Now we can write $p = \frac{1}{2}q + \frac{1}{2}r$, where

$$q(x) = \begin{cases} p(x), & x \neq x', x \neq x'' \\ p(x) + \varepsilon, & x = x' \\ p(x) - \varepsilon, & x = x'' \end{cases} \quad (270)$$

$$r(x) = \begin{cases} p(x), & x \neq x', x \neq x'' \\ p(x) - \varepsilon, & x = x' \\ p(x) + \varepsilon, & x = x'' \end{cases} \quad (271)$$

and

$$\varepsilon = \frac{1}{2} \min \left\{ p(x'), p(x''), \frac{1}{n+1} - p(x'), \frac{1}{n+1} - p(x'') \right\}. \quad (272)$$

Thus, $p \notin \text{extr}(\mathcal{A}_n(\bar{x}))$, giving the contradiction. We have shown $\mathcal{F} = \text{extr}(\mathcal{A}_n(\bar{x}))$ and part 1 of the lemma.

To prove part 2 of the lemma, first define

$$\mathcal{B}_n \triangleq \bigcup_{x \in \mathcal{X}} \mathcal{F}_n(x). \quad (273)$$

For any n , it holds that

$$\mathcal{B}_{n+1} \subseteq \text{conv}(\mathcal{B}_n). \quad (274)$$

This follows from writing $p \in \mathcal{B}_{n+1}$ as

$$p = \sum_{q \in \mathcal{B}_n : \text{supp}(q) \subseteq \text{supp}(p)} \frac{1}{n+1} q. \quad (275)$$

One can establish part 2 by using part 1 and (274). \blacksquare

B. Proof of Theorem 2

With Lemma 8 in hand, we are equipped to prove Theorem 2. Fix R_0 and let U^* be the maximizer of $\Pi(R_0)$. When the payoff function is $\pi(x, z) = \mathbf{1}\{x \neq z\}$, we can rewrite $\Pi(R_0)$ as

$$\Pi(R_0) = \min_{z(u)} \mathbb{E} \pi(X, z(U^*)) \quad (276)$$

$$= \min_{z(u)} \sum_u P_{U^*}(u) \sum_x P_{X|U^*}(x|u) \mathbf{1}\{x \neq z(u)\} \quad (277)$$

$$= \sum_u P_{U^*}(u) \min_z \sum_x P_{X|U^*}(x|u) \mathbf{1}\{x \neq z\} \quad (278)$$

$$= \sum_u P_{U^*}(u) \min_z (1 - P_{X|U^*}(z|u)) \quad (279)$$

$$= \sum_u P_{U^*}(u) (1 - \max_x P_{X|U^*}(x|u)). \quad (280)$$

We now show that the set $\{P_{X|U^*=u}\}_u$ in (38) can be restricted to the finite set $\mathcal{P}_{\text{unif}}$, where

$$\mathcal{P}_{\text{unif}} \triangleq \{p \in \Delta_{\mathcal{X}} : p = \text{Unif}(\mathcal{A}) \text{ for some } \mathcal{A} \subseteq \mathcal{X}\}. \quad (281)$$

By applying part 2 of Lemma 8 to each distribution in $\{P_{X|U^*=u}\}_u$, we have that there exists a random variable V such that

$$\forall u, v, P_{X|U^*=u, V=v} \in \mathcal{P}_{\text{unif}} \quad (282)$$

$$\begin{aligned} \forall u, v, v', \arg \max_x P_{X|U^*V}(x|u, v) \\ = \arg \max_x P_{X|U^*V}(x|u, v'). \end{aligned} \quad (283)$$

We now write

$$\begin{aligned} \Pi(R_0) \\ \stackrel{(a)}{=} \sum_u P_{U^*}(u) (1 - \max_x P_{X|U^*}(x|u)) \end{aligned} \quad (284)$$

$$= \sum_u P_{U^*}(u) (1 - \max_x \sum_v P_{X|U^*V}(x|u, v) P_{V|U^*}(v|u)) \quad (285)$$

$$\stackrel{(b)}{=} \sum_{u,v} P_{U^*V}(u, v) (1 - \max_x P_{X|U^*V}(x|u, v)) \quad (286)$$

$$= \min_{z(u,v)} \mathbb{E} \pi(X, z(U^*, V)), \quad (287)$$

where (a) is due to (280) and (b) follows from (283). By noting that $R_0 \geq H(X|U^*) \geq H(X|U^*, V)$ and letting $U = (U^*, V)$, we have

$$\Pi(R_0) \leq \max_{\substack{U: P_{X|U^*=u} \in \mathcal{P}_{\text{unif}} \\ R_0 \geq H(X|U)}} \min_{z(u)} \mathbb{E} \pi(X, z(U)). \quad (288)$$

This shows that we can restrict attention to $\mathcal{P}_{\text{unif}}$ without hurting the payoff. Now, observe that $p \in \mathcal{P}_{\text{unif}}$ satisfies

$$(H(p), 1 - \max_x p(x)) = (\log n, \frac{n-1}{n}) \quad (289)$$

for some $n \in \mathbb{N}$. Referring to (280) and noting that $H(X|U) = \sum_u P_U(u) H(X|U=u)$, we see that $\Pi(R_0)$ cannot lie outside of the convex hull of the pairs $(\log n, \frac{n-1}{n}), n \in \mathbb{N}$. That is,

$$\Pi(R_0) \leq \phi(R_0). \quad (290)$$

To see $\Pi(R_0) \leq \pi_{\text{max}}$, simply write

$$\Pi(R_0) = \sum_u P_{U^*}(u) (1 - \max_x P_{X|U^*}(x|u)) \quad (291)$$

$$\leq 1 - \max_x \sum_u P_{U^*}(u) P_{X|U^*}(x|u) \quad (292)$$

$$= \pi_{\text{max}}. \quad (293)$$

It remains to show that $\min\{\phi(R_0), \pi_{\text{max}}\}$ can be achieved through the proper choice of U . To that end, let \bar{x} and N be such that $P_X \in \mathcal{A}_N(\bar{x})$. By the convexity of \mathcal{R} , we will be done once we show that we can achieve not only the points $(\log n, \frac{n-1}{n}), n \in [N]$, but also the intersection of ϕ with π_{max} . To achieve the point $(\log n, \frac{n-1}{n})$, invoke part 2 of Lemma 8 produce U . Denote the corresponding rate-payoff pair by (R'_0, Π') . Since the $\{P_{X|U=u}\}_u$ all satisfy

$$(H(X|U=u), 1 - \max_x P_{X|U=u}(x|u)) = (\log n, \frac{n-1}{n}) \quad (294)$$

so must (R'_0, Π') as well. To achieve the intersection of ϕ with π_{max} , first invoke part 1 of Lemma 8 to produce U . Denote the corresponding rate-payoff pair by (R''_0, Π'') . The $\{P_{X|U=u}\}_u$ correspond to either $(\log n, \frac{n-1}{n})$ or $(\log(n+1), \frac{n}{n+1})$. Thus, (R''_0, Π'') lies on f because it is a convex combination of those two points. We also have that (R''_0, Π'') satisfies $\Pi'' = \pi_{\text{max}}$ because

$$\arg \max_x P_{X|U=u}(x|u) = \bar{x}, \forall u \in \mathcal{U}. \quad (295)$$

This completes the proof of Theorem 2.

APPENDIX B PROOF OF LEMMA 5

Let $R_0 > I(W; V|U)$. Define the typical set

$$\mathcal{T}_\varepsilon^n \triangleq \{u^n : |T_{u^n}(u) - P_U(u)| < \varepsilon P_U(u), \forall u \in \mathcal{U}\}. \quad (296)$$

where T_{u^n} denotes the type of u^n .

First, write

$$\begin{aligned} & \left\| Q_{MW^n X_B Y_B} - \widehat{Q}_{MW^n X_B Y_B} \right\| \\ &= \sum_{m: U^n(m) \in \mathcal{T}_\varepsilon^n} Q_M(m) \left\| Q_{W^n X_B Y_B | M=m} - \widehat{Q}_{W^n X_B Y_B | M=m} \right\| \\ & \quad + \sum_{m: U^n(m) \notin \mathcal{T}_\varepsilon^n} Q_M(m) \left\| Q_{W^n X_B Y_B | M=m} - \widehat{Q}_{W^n X_B Y_B | M=m} \right\|. \end{aligned} \quad (297)$$

The expected value of the second term in (297) can be bounded easily. For sufficiently large n , we have

$$\begin{aligned} & \mathbb{E} \sum_{m: U^n(m) \notin \mathcal{T}_\varepsilon^n} Q_M(m) \left\| P_{X^n Y^k | M=m} - \widehat{Q}_{X^n Y^k | M=m} \right\| \\ & \leq \mathbb{E} \sum_{m: U^n(m) \notin \mathcal{T}_\varepsilon^n} Q_M(m) \end{aligned} \quad (298)$$

$$= \mathbb{P}[U^n(M) \notin \mathcal{T}_\varepsilon^n] \quad (299)$$

$$= \mathbb{P}[U^n(1) \notin \mathcal{T}_\varepsilon^n] \quad (300)$$

$$\stackrel{(a)}{\leq} \varepsilon, \quad (301)$$

where (a) is due to the law of large numbers.

The expected value of the first term in (297) can first be rewritten by moving the expectation with respect to the subcodebook $C_V^{(n)}(m)$ inside the sum.

$$\begin{aligned} & \mathbb{E} \sum_{m: U^n(m) \in \mathcal{T}_\varepsilon^n} Q_M(m) \left\| Q_{W^n X_B Y_B | M=m} - \widehat{Q}_{W^n X_B Y_B | M=m} \right\| \\ &= \mathbb{E}_{C_V^{(n)}} \sum_{m: U^n(m) \in \mathcal{T}_\varepsilon^n} Q_M(m) \mathbb{E}_{C_V^{(n)}(m)} \left\| Q_{W^n X_B Y_B | M=m} \right. \\ & \quad \left. - \widehat{Q}_{W^n X_B Y_B | M=m} \right\|. \end{aligned} \quad (302)$$

It remains to show that the inner expectation vanishes for each m .¹¹

To do this, first observe that $Q_{W^n X_B Y_B | M=m}$ is the output of the memoryless (but nonstationary) channel $\Phi \triangleq Q_{W^n X_B Y_B | K, M=m}$ acting on a codebook of size 2^{nR_0} that

¹¹Due to the symmetry of codebook construction, the behavior of the inner expectation is uniform for all m . Thus, the rate of convergence does not play a role in claiming that (302) vanishes.

is generated i.i.d. according to $\Psi \triangleq \prod_i P_{V|U=u_i(m)}$. Furthermore, it can be verified that $\widehat{Q}_{W^n X_B Y_B | M=m}$ is the output distribution of the channel Φ when the input distribution is Ψ . Thus, we can invoke the soft covering lemma (Lemma 3) as long as R_0 exceeds the sup-information rate of the process that results from Φ acting on Ψ . To be explicit, that process is given by

$$\begin{aligned} & \Gamma(v^n, w^n, x_B, y_B) \\ & \triangleq \prod_{i=1}^n P_{VW|U}(w_i, v_i | u_i(m)) \prod_{i \in \mathcal{B}} P_{VXY|U}(x_i, y_i, v_i | u_i(m)). \end{aligned} \quad (303)$$

Since Γ is a memoryless process and the second moments of $\{i_\Gamma(W_i, X_i, Y_i; V_i)\}$ are uniformly bounded, the law of large numbers gives

$$\limsup_{n \rightarrow \infty} \frac{1}{n} i_\Gamma(W^n, X_B, Y_B; V^n) \leq \mathbb{E} \frac{1}{n} i_\Gamma(W^n, X_B, Y_B; V^n). \quad (304)$$

Furthermore, we can upper bound the expected information density by writing

$$\begin{aligned} & \mathbb{E} \frac{1}{n} i_\Gamma(W^n, X_B, Y_B; V^n) \\ & = \mathbb{E} \frac{1}{n} i_\Gamma(W^n; V^n) + \mathbb{E} \frac{1}{n} i_\Gamma(X_B, Y_B; V^n | W^n) \end{aligned} \quad (305)$$

$$= \mathbb{E} \frac{1}{n} i_\Gamma(W^n; V^n) + \frac{1}{n} I_\Gamma(X_B, Y_B; V^n | W^n) \quad (306)$$

$$\leq \mathbb{E} \frac{1}{n} i_\Gamma(W^n; V^n) + \alpha \log |\mathcal{X}| |\mathcal{Y}| \quad (307)$$

$$\begin{aligned} & = \mathbb{E} \frac{1}{n} \sum_{i=1}^n i_{P_{WV|U=u_i(m)}}(W; V | U = u_i(m)) \\ & \quad + \alpha \log |\mathcal{X}| |\mathcal{Y}| \end{aligned} \quad (308)$$

$$= \frac{1}{n} \sum_{i=1}^n I(W; V | U = u_i(m)) + \alpha \log |\mathcal{X}| |\mathcal{Y}| \quad (309)$$

$$= \sum_{u \in \mathcal{U}} T_{u^n(m)} I(W; V | U = u) + \alpha \log |\mathcal{X}| |\mathcal{Y}| \quad (310)$$

$$\stackrel{(a)}{\leq} \sum_{u \in \mathcal{U}} (1 + \varepsilon) P_U(u) I(W; V | U = u) + \alpha \log |\mathcal{X}| |\mathcal{Y}| \quad (311)$$

$$= (1 + \varepsilon) I(W; V | U) + \alpha \log |\mathcal{X}| |\mathcal{Y}|. \quad (312)$$

Step (a) follows from $u^n(m) \in \mathcal{T}_\varepsilon^n$.

The expression in (312) is strictly less than R_0 for the proper choice of $\varepsilon > 0$ and $\alpha > 0$. Thus, when $u^n(m) \in \mathcal{T}_\varepsilon^n$,

$$R_0 > \limsup_{n \rightarrow \infty} \frac{1}{n} i_\Gamma(W^n, X_B, Y_B; V^n). \quad (313)$$

Invoking Lemma 3, we have

$$\lim_{n \rightarrow \infty} \mathbb{E}_{\mathcal{C}_V^{(n)}(m)} \left\| Q_{W^n X_B Y_B | M=m} - \widehat{Q}_{W^n X_B Y_B | M=m} \right\| = 0. \quad (314)$$

This completes the proof of Lemma 5.

REFERENCES

[1] P. Cuff, "A framework for partial secrecy," in *Proc. Global Telecomm. Conf. (GLOBECOM)*, Dec. 2010.

[2] —, "Using a secret key to foil an eavesdropper," in *48th Allerton Conf. on Communication, Control, and Computing*, Sept. 2010, pp. 1405–1411.

[3] C. Schieler and P. Cuff, "Secrecy is cheap if the adversary must reconstruct," in *Proc. IEEE International Symposium on Information Theory (ISIT)*, Jul. 2012, pp. 66–70.

[4] P. Cuff, "Optimal equivocation in secrecy systems a special case of distortion-based characterization," in *Information Theory and Applications Workshop (ITA)*, Feb. 2013, pp. 1–3.

[5] C. Schieler and P. Cuff, "Rate-distortion theory for secrecy systems," in *Proc. IEEE International Symposium on Information Theory (ISIT)*, Jul. 2013, pp. 2219–2223.

[6] C. Shannon, "Communication theory of secrecy systems," *Bell Syst. Tech. J.*, vol. 28, no. 4, pp. 656–715, Oct. 1949.

[7] A. Wyner, "The wire-tap channel," *Bell Syst. Tech. J.*, vol. 54, no. 8, pp. 1334–1387, 1975.

[8] I. Csiszár and J. Körner, "Broadcast channels with confidential messages," *IEEE Trans. Inf. Theory*, vol. 24, no. 3, pp. 339–348, May 1978.

[9] U. Maurer, "Secret key agreement by public discussion from common information," *IEEE Trans. Inf. Theory*, vol. 39, no. 3, pp. 733–742, May 1993.

[10] R. Ahlswede and I. Csiszár, "Common randomness in information theory and cryptography. i. secret sharing," *IEEE Trans. Inf. Theory*, vol. 39, no. 4, pp. 1121–1132, Jul. 1993.

[11] H. Yamamoto, "A rate-distortion problem for a communication system with a secondary decoder to be hindered," *IEEE Trans. Inf. Theory*, vol. 34, no. 4, pp. 835–842, July 1988.

[12] —, "Rate-distortion theory for the Shannon cipher system," *IEEE Trans. Inf. Theory*, vol. 43, no. 3, pp. 827–835, May 1997.

[13] P. Cuff, H. Permuter, and T. Cover, "Coordination capacity," *IEEE Trans. Inf. Theory*, vol. 56, no. 9, pp. 4181–4206, Sept. 2010.

[14] U. Maurer, A. Ruedlinger, and B. Tackmann, "Confidentiality and integrity: a constructive perspective," in *Theory of Cryptography*, ser. Lecture Notes in Computer Science, R. Cramer, Ed., 2012, vol. 7194, pp. 209–229.

[15] T. Courtade and T. Weissman, "Multiterminal source coding under logarithmic loss," *IEEE Trans. Inf. Theory*, 2013, to appear.

[16] P. Cuff, "Distributed channel synthesis," *IEEE Trans. Inf. Theory*, vol. 59, no. 11, pp. 7071–7096, Nov. 2013.

[17] A. Wyner, "The common information of two dependent random variables," *IEEE Trans. Inf. Theory*, vol. 21, no. 2, pp. 163–179, Mar. 1975.

[18] T. S. Han and S. Verdú, "Approximation theory of output statistics," *IEEE Trans. Inf. Theory*, vol. 39, no. 3, pp. 752–772, May 1993.

[19] M. Yassaee, M. Aref, and A. Gohari, "Achievability proof via output statistics of random binning," *IEEE Trans. Inf. Theory*, vol. PP, no. 99, pp. 1–1, 2014.

[20] A. Winter, "Secret, public and quantum correlation cost of triples of random variables," in *Proc. IEEE International Symposium on Information Theory (ISIT)*, Sept. 2005, pp. 2270–2274.

[21] I. Csiszár and J. Körner, *Information theory: coding theorems for discrete memoryless systems*. Cambridge University Press, 2011.