

On Random Coding Error Exponents of Watermarking Codes

Neri Merhav*
Hewlett-Packard Laboratories Israel
HPL-98-169
October, 1998

E-mail: merhav@hp.technion.ac.il

steganography,
watermarking,
information hiding,
error exponent,
random coding

Watermarking codes are analyzed from an information-theoretic viewpoint as a game between an information hider and an active attacker. While the information hider embeds a secret message (watermark) in a covert message (typically, text, image, sound, or video stream) within a certain distortion level, the attacker processes the resulting watermarked message, within limited additional distortion, in an attempt to invalidate the watermark. For a memoryless covert source, we provide a single-letter characterization of the minimax-maximin game of the random coding error exponent associated with the average probability of erroneously decoding the watermark. This single-letter characterization is in effect because there is a “memoryless saddle point” in this game: The information hider utilizes a memoryless channel to generate random codewords for every covert message, whereas the attacker implements a memoryless channel to disrupt the watermark information hidden in the covert message.

Internal Accession Date Only

*Electrical Engineering Department of the Technion–Israel Institute of Technology, Technion City, Haifa 32000 Israel

Part of this work was done while the author was a consultant at Hewlett-Packard Laboratories, Palo Alto, California

© Copyright Hewlett-Packard Company 1998

1 Introduction

The fast development of digital information technology, combined with the simplicity of duplication and distribution of digital data, has recently stimulated many research efforts towards the design and study of sophisticated copyright protection and information hiding methodologies (see, e.g., [1],[2],[4] and references therein).

Relatively little attention, however, has been devoted to the problem of information hiding from an information-theoretic perspective. An exception is a recent work by O'Sullivan, Moulin, and Ettinger [7], who characterized the highest achievable information rate of watermarking codes for the following system model: A secret message, encoded at rate R , is hidden in a memoryless covertext message within small degradation of quality, symbolized by distortion level D_1 w.r.t. some distortion measure. An active attacker may introduce additional distortion D_2 in attempt to disrupt the watermark. Finally, the resulting data set is analyzed using information shared with the information hider (i.e., the original covertext message) to extract the watermark. In [7], the highest information rate R has been found (as a function of D_1 and D_2) for which, even in the presence of optimum attack, there exist block encoders and decoders that guarantee arbitrarily small probability of erroneous watermark decoding.

In this paper, we go one step further and provide a single-letter characterization of the best achievable random coding error exponent under this model. Our results are based heavily on a fact, proved in Section 3, that for a memoryless covertext source and a memoryless random coding distribution, an optimum attack strategy subject to the above-mentioned distortion constraint, is to implement a *memoryless channel* on the watermarked data set. It is then clear that this memoryless channel is the one that minimizes the random coding error exponent subject to the single-letter

distortion constraint. We also show in Section 4 the dual property that if the attack strategy is given by a memoryless channel, then the optimal random coding distribution is memoryless as well. It turns out then that there is a minimax–maximin saddle point in the game between the information hider and the attacker, where in the presence of a memoryless covertext source, both parties implement memoryless channels.

2 Notation, Problem Formulation, and Preliminaries

Throughout this paper, we adopt the following notation conventions.

Scalar random variables will be denoted by capital letters (e.g., X) and specific values they may take will be denoted by the corresponding lower case letters (e.g., x). All scalar random variables in this paper are assumed to take on values in the same finite alphabet \mathcal{A} . Similarly, random vectors of length n (n – positive integer) will be denoted by boldface capital letters (e.g., $\mathbf{X} = (X_1, \dots, X_n)$), and specific values that they take will be denoted by the respective boldface lower case letters (e.g., $\mathbf{x} = (x_1, \dots, x_n)$), where it is understood that the alphabet of each such vector is always \mathcal{A}^n , the n th order Cartesian power of \mathcal{A} .

Probability mass functions (PMF’s) of single letters will also be denoted by capital letters with probabilities of specific letters denoted by the respective lower case letters (e.g., $Q = \{q(x), x \in \mathcal{A}\}$). Similarly, vector sources, or joint PMF’s of n -vectors, will be denoted by boldface capital letters with probabilities of specific vector values denoted by the corresponding boldface lower case letters (e.g., $\mathbf{Q} = \{\mathbf{q}(\mathbf{x}), \mathbf{x} \in \mathcal{A}^n\}$). A source \mathbf{Q} over \mathcal{A}^n is said to assume a *product form* if $\mathbf{Q}(\mathbf{x}) = \prod_{i=1}^n q(x_i)$ for all $\mathbf{x} \in \mathcal{A}^n$. For shorthand notation, this fact will be denoted by $\mathbf{Q} = Q^n$, where $Q = \{q(x), x \in \mathcal{A}\}$ will be referred to as the *single-letter component* of \mathbf{Q} .

Similar conventions apply to conditional PMF's of channels: Single-letter conditional PMF's will be denoted by capital letters that symbolize matrices of conditional probabilities (e.g., $W = \{w(y|x), x, y \in \mathcal{A}\}$) and vector channels will be denoted by the respective boldface letters (e.g., $\mathbf{W} = \{\mathbf{w}(\mathbf{y}|\mathbf{x}), \mathbf{x}, \mathbf{y} \in \mathcal{A}^n\}$). A channel \mathbf{W} is said to admit a product form if $\mathbf{w}(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^n w(y_i|x_i)$ for all $\mathbf{x}, \mathbf{y} \in \mathcal{A}^n$. For shorthand notation, this fact will be denoted by $\mathbf{W} = W^n$, where $W = \{w(y|x), x, y \in \mathcal{A}\}$ will be referred to as the single-letter component of \mathbf{W} .

The probability of an event \mathcal{E} will be denoted by $\Pr\{\mathcal{E}\}$ whenever the underlying probability measure is clear from the context, and the expectation operator will be denoted by $\mathbf{E}\{\cdot\}$. Summations and products of indexed terms over the entire index set will be sometimes subscripted only by the symbol of that index. For example, $\sum_x f(x)$ will mean that the summation is over the entire alphabet \mathcal{A} , $\prod_i q(x_i)$ means that the product is taken from $i = 1$ to n , etc.

Let $\mathbf{P} = P^n$ designate a stationary, memoryless vector source generating random vectors $\mathbf{U} = (U_1, \dots, U_n)$ over \mathcal{A}^n . The data vector \mathbf{U} will designate the *covert* message within which the watermark will be hidden. Let $d_1 : \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}_+$ denote a single-letter distortion measure.

A rate R *watermarking code* of size n , subject to distortion level D_1 , is a triple $(\mathcal{M}_n, f_n, g_n)$ with the following ingredients. The set \mathcal{M}_n is the watermarking message set whose size is $\lfloor 2^{nR} \rfloor$. The function $f_n : \mathcal{A}^n \times \mathcal{M}_n \rightarrow \mathcal{A}^n$ is the encoder that maps every combination of a covert data set $\mathbf{u} \in \mathcal{A}^n$ and a *watermark message* $m \in \mathcal{M}_n$ into a *stegotext* message (or, watermarked message) $\mathbf{x} \in \mathcal{A}^n$ such that $\mathbf{E}d_1(u_i, X_i) \leq D_1$ for all $\mathbf{u} \in \mathcal{A}^n$, $1 \leq i \leq n$, that is, the conditional expected distortion in each coordinate given $\mathbf{U} = \mathbf{u}$ never exceeds D_1 for all \mathbf{u} .¹ The function $g_n : \mathcal{A}^n \times \mathcal{A}^n \rightarrow \mathcal{M}_n$ is the decoder, which maps the original covert message (shared with the

¹Although these constraints seem quite restrictive, note that for stationary additive channels and difference distortion measures, $d_1(u, x) = d(x - u)$, they all boil down to a single constraint $\mathbf{E}d(V) \leq D_1$, where $V = X - U$ designates the additive noise.

encoder) together with a stegotext message (possibly modified or corrupted by an attacker) back to a watermark message.

Let $d_2 : \mathcal{A} \times \mathcal{A} \rightarrow \mathbb{R}_+$ denote a single-letter distortion measure. An *attack* subject to distortion D_2 is a (possibly randomized) map $h_n : \mathcal{A}^n \rightarrow \mathcal{A}^n$ that satisfies $\mathbf{E}d_2(x_i, Y_i) \leq D_2$ for every $\mathbf{x} \in \mathcal{A}^n$, $1 \leq i \leq n$. We will also assume throughout the paper that $\max\{\max_{u,x} d_1(u, x), \max_{x,y} d_2(x, y)\} \triangleq D_{\max} < \infty$.

Given distortion levels D_1 and D_2 w.r.t. distortion measures d_1 and d_2 , respectively, the error probability is defined as

$$P_e(f_n, g_n, h_n) = \Pr\{g_n(\mathbf{U}, h_n(f_n(\mathbf{U}, M))) \neq M\}, \quad (1)$$

where M designates a random variable in \mathcal{M}_n with a uniform PMF independently of \mathbf{U} .

An *achievable rate* R is one for which there exists a sequence of rate R watermarking codes $\{(\mathcal{M}_n, f_n, g_n)\}_{n \geq 1}$ subject to distortion level D_1 such that

$$\lim_{n \rightarrow \infty} \sup_{h_n} P_e(f_n, g_n, h_n) = 0, \quad (2)$$

where the supremum is over all attack strategies subject to distortion D_2 . The capacity of the system, $C = C(D_1, D_2)$, is the supremum of all achievable rates.

Similarly as in [7],² it can be shown that for a memoryless covertext source $\mathbf{P} = P^n$, the capacity is given by

$$C(D_1, D_2) = \max \min I(X; Y|U) \quad (3)$$

where the maximin is taken over all Markov chains $U \ominus X \ominus Y$ such that: (i) the marginal PMF of U is given by P , (ii) $\mathbf{E}d_1(u, X) \leq D_1$ for all $u \in \mathcal{A}$, and (iii) $\mathbf{E}d_2(x, Y) \leq D_2$ for all $x \in \mathcal{A}$.

²The above-described model assumptions are somewhat different from those of [7].

In this paper, we refine the analysis of the probability of error at a given information rate $R < C(D_1, D_2)$. To this end, we assume a random coding regime and we focus on the average probability of error w.r.t. the ensemble of randomly chosen codes given \mathbf{U} . To meet the D_1 distortion constraint, we adopt the following random coding mechanism: For each $\mathbf{u} \in \mathcal{A}^n$, we generate independently, at random 2^{nR} codewords by using a certain channel $\mathbf{Q} = \{\mathbf{q}(\mathbf{x}|\mathbf{u}), \mathbf{u}, \mathbf{x} \in \mathcal{A}^n\} \in \mathcal{Q}_n(D_1)$ fed by \mathbf{u} , where

$$\mathcal{Q}_n(D_1) \triangleq \{\mathbf{Q} : \sum_{\mathbf{x}} \mathbf{q}(\mathbf{x}|\mathbf{u}) d_1(u_i, x_i) \leq D_1 \quad \forall \mathbf{u} \in \mathcal{A}^n, 1 \leq i \leq n\}.$$

The channel \mathbf{Q} will be referred to as the *watermarking channel* [7]. For $n = 1$, $\mathcal{Q}_n(D_1) = \mathcal{Q}_1(D_1)$, which is a set of single-letter channels Q , will be abbreviated by $\mathcal{Q}(D_1)$.

Since the attacker is allowed to use a randomized map, we will symbolize the attack strategy by a *channel* $\mathbf{W} = \{\mathbf{w}(\mathbf{y}|\mathbf{x}), \mathbf{x}, \mathbf{y} \in \mathcal{A}^n\}$, henceforth referred to as the *attack channel* [7]. The fact that the attack is subjected to distortion constraint D_2 means that \mathbf{W} lies in the set

$$\mathcal{W}_n(D_2) = \{\mathbf{W} : \sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) d_2(x_i, y_i) \leq D_2 \quad \forall \mathbf{x} \in \mathcal{A}^n, 1 \leq i \leq n\}.$$

For $n = 1$, $\mathcal{W}_n(D_2) = \mathcal{W}_1(D_2)$, which is a set of single-letter channels W , will be abbreviated by $\mathcal{W}(D_2)$.

The average error probability $\bar{P}_e(\mathbf{Q}, \mathbf{W})$ is defined as the expectation of P_e w.r.t. the joint ensemble of random codes drawn by \mathbf{Q} and the randomized attacks governed by \mathbf{W} .

We will be interested in the quantities

$$\min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \bar{P}_e(\mathbf{Q}, \mathbf{W})$$

and

$$\max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \bar{P}_e(\mathbf{Q}, \mathbf{W})$$

as well as on a characterization of pairs (\mathbf{Q}, \mathbf{W}) that achieve these minimax and maximin criteria. However, since an exact expression of $\bar{P}_e(\mathbf{Q}, \mathbf{W})$ is unavailable, we will focus on a Gallager-type upper bound [6] on $\bar{P}_e(\mathbf{Q}, \mathbf{W})$.

3 The Worst Attack Channel is Memoryless

In this section, we confine attention to the minimax criterion, corresponding to the assumption that the attacker has perfect knowledge of the watermarking channel, which is assumed here to admit a product form $\mathbf{Q} = Q^n$. Generally speaking, a more plausible assumption would be that the attacker knows not only \mathbf{Q} , but also the specific watermark code in use. Our above assumption on knowing Q but not the code itself is suitable, for example, if the random choice of the code is controlled by a secret key shared by the encoder and decoder, but not the attacker.

The first observation to make is that the problem at hand is no different from an ordinary channel coding problem with side information at both encoder and decoder: A message $m \in \mathcal{M}_n$ of length nR bits is mapped to a codeword $\mathbf{x} \in \mathcal{A}^n$ depending on the side information $\mathbf{u} \in \mathcal{A}^n$. (In other words, there is a different codebook for every \mathbf{u} .) Then, \mathbf{x} is transmitted across a channel \mathbf{W} , whose output \mathbf{y} is processed in the presence of \mathbf{u} to decode the message.

Assuming the random coding regime described in Section 2 and maximum likelihood decoding, the following upper bound on the average error probability is easily derived as a straightforward extension of [6, Theorem 5.6.1, p. 135] to the case where side information is available at both encoder and decoder:

$$\bar{P}_e(\mathbf{Q}, \mathbf{W}) \leq B(\rho, \mathbf{Q}, \mathbf{W}) \triangleq 2^{\rho n R} F(\rho, \mathbf{Q}, \mathbf{W}), \quad 0 \leq \rho \leq 1 \quad (4)$$

where

$$F(\rho, \mathbf{Q}, \mathbf{W}) \triangleq \sum_{\mathbf{u}} \mathbf{p}(\mathbf{u}) \sum_{\mathbf{y}} \left[\sum_{\mathbf{x}} q(\mathbf{x}|\mathbf{u}) w(\mathbf{y}|\mathbf{x})^{1/(1+\rho)} \right]^{1+\rho} \quad (5)$$

with $\mathbf{P} = \{\mathbf{p}(\mathbf{u}), \mathbf{u} \in \mathcal{A}^n\}$ being a given memoryless source P^n , $\mathbf{Q} \in \mathcal{Q}_n(D_2)$, and $\mathbf{W} \in \mathcal{W}_n(D_2)$.

The objective function hereafter will be the upper bound on the average error probability $B(\rho, \mathbf{Q}, \mathbf{W})$. This function involves the auxiliary parameter ρ , which can be chosen so as to obtain the tightest upper bound, namely, $\min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W})$. Thus, for a memoryless watermarking channel $\mathbf{Q} = Q^n$, we will be interested in characterizing the minimax game

$$\min_{\mathbf{Q} \in \mathcal{Q}(D_1)} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{0 \leq \rho \leq 1} B(\rho, Q^n, \mathbf{W}). \quad (6)$$

For later use, we also define $F(\rho, Q, W)$ and $B(\rho, Q, W)$ as the same functionals as above, applied to the case $n = 1$, that is:

$$B(\rho, Q, W) = 2^{\rho R} F(\rho, Q, W) = 2^{\rho R} \sum_u p(u) \sum_y \left[\sum_x q(x|u) w(y|x)^{1/(1+\rho)} \right]^{1+\rho}. \quad (7)$$

Note that for $\mathbf{P} = P^n$, $B(\rho, Q^n, W^n) = B(\rho, Q, W)^n$ and $F(\rho, Q^n, W^n) = F(\rho, Q, W)^n$. We also define the random coding error exponent as

$$E_r(\rho, Q, W) \triangleq -\log B(\rho, Q, W) = -\log F(\rho, Q, W) - \rho R. \quad (8)$$

Our first result states that for a memoryless covertext source and a memoryless random coding distribution, the worst channel \mathbf{W} , in the sense of maximizing B subject a distortion constraint, is memoryless as well. Since B and F are the same within a factor that does depends on \mathbf{Q} and \mathbf{W} , we can state this result in terms of F .

Theorem 1 *Given ρ , $\mathbf{P} = P^n$, Q , and distortion level D_2 w.r.t. distortion measure d_2 , it is a stationary memoryless channel $\mathbf{W}_* = W_*^n$ that attains the maximum of $F(\rho, Q^n, \mathbf{W})$ w.r.t. \mathbf{W}*

subject to the constraints

$$\mathbf{w}(\mathbf{y}|\mathbf{x}) \geq 0 \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{A}^n \quad (9)$$

$$\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) = 1 \quad \forall \mathbf{x} \in \mathcal{A}^n \quad (10)$$

$$\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) d_2(x_i, y_i) \leq D_2, \quad \forall \mathbf{x} \in \mathcal{A}^n, 1 \leq i \leq n. \quad (11)$$

The single-letter component W_* of \mathbf{W}_* is given by the maximizer of $F(\rho, Q, W)$ subject to the same constraints applied to the case $n = 1$, i.e.,

$$\mathbf{W} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} F(\rho, Q^n, \mathbf{W}) = \left[\max_{W \in \mathcal{W}(D_2)} F(\rho, Q, W) \right]^n. \quad (12)$$

The proof is given in Section 5. Theorem 1 leads to our first main result, which is the following.

Theorem 2 For a memoryless source $\mathbf{P} = P^n$,

$$\min_{Q \in \mathcal{Q}(D_2)} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{0 \leq \rho \leq 1} B(\rho, Q^n, \mathbf{W}) = 2^{-nE(R, D_1, D_2)}, \quad (13)$$

where

$$E(R, D_1, D_2) = \max_{0 \leq \rho \leq 1} \max_{Q \in \mathcal{Q}(D_1)} \min_{W \in \mathcal{W}(D_2)} E_r(\rho, Q, W). \quad (14)$$

Note that the order of the optimizations in eq. (14) is different from the one in eq. (13). The order in (14) has the convenience that for every fixed ρ , the functional $B(\rho, Q, W) = 2^{-E_r(\rho, Q, W)}$ as we shall see, is concave in W and convex in Q .

Proof of Theorem 2. For a given $Q \in \mathcal{Q}(D_2)$, consider the expression

$$\max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{0 \leq \rho \leq 1} B(\rho, Q^n, \mathbf{W}) = \exp_2 \left\{ \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{0 \leq \rho \leq 1} [\log F(\rho, Q^n, \mathbf{W}) + \rho n R] \right\}. \quad (15)$$

Since $F(\rho, Q^n, \mathbf{W})$ is concave in \mathbf{W} (see Section 5, Lemma 1) and the logarithmic function is monotonic and concave, then $\log F(\rho, Q^n, \mathbf{W})$ is concave as well. Also, a straightforward extension

of [5, p. 17, Proof of Theorem 2] to the case of side information, shows that $\log F(\rho, Q^n, \mathbf{W})$, and hence also $[\log F(\rho, Q^n, \mathbf{W}) + \rho n R]$ is convex in ρ . Also, $\mathcal{W}_n(D_2)$ and the interval $[0, 1]$ are both compact convex sets. Therefore, by the minimax theorem [8, p. 232, Theorem (6.3.7)], the maximization and the minimization on the r.h.s. of eq. (15) are interchangeable, and so, it can be written also as

$$\begin{aligned} \exp_2 \left\{ \min_{0 \leq \rho \leq 1} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} [\log F(\rho, Q^n, \mathbf{W}) + \rho n R] \right\} &= \min_{0 \leq \rho \leq 1} \left[2^{\rho n R} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} F(\rho, Q^n, \mathbf{W}) \right] \\ &= \left[\min_{0 \leq \rho \leq 1} 2^{\rho R} \max_{W \in \mathcal{W}(D_2)} F(\rho, Q, W) \right]^n \quad (16) \end{aligned}$$

where the second equality follows from Theorem 1. The assertion of the theorem now follows by taking the minimum over $Q \in \mathcal{Q}(D_2)$. \square

4 The Best Watermarking Channel is Memoryless

In this section, we focus on the maximin criterion, which corresponds to a situation where the information hider knows in advance the attack channel and strives at maximizing the error exponent for this channel. In analogy to the previous section, here we assume that the attack channel is memoryless, and so, our objective function will be

$$\max_{W \in \mathcal{W}(D_2)} \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, W^n).$$

We next show that the solution to this problem is dual to the solution of minimax problem of Section 3. As a first step towards this end, we need the following analogue to Theorem 1.

Theorem 3 *Given ρ , $\mathbf{P} = P^n$, W , and a distortion level D_1 w.r.t. distortion measure d_1 , it is a stationary memoryless channel $\mathbf{Q}_* = Q_*^n$ that attains the minimum of $F(\rho, \mathbf{Q}, W^n)$ w.r.t. \mathbf{Q} subject*

to the constraints

$$\mathbf{q}(\mathbf{x}|\mathbf{u}) \geq 0 \quad \forall \mathbf{u}, \mathbf{x} \in \mathcal{A}^n \quad (17)$$

$$\sum_{\mathbf{x}} \mathbf{q}(\mathbf{x}|\mathbf{u}) = 1 \quad \forall \mathbf{u} \in \mathcal{A}^n \quad (18)$$

$$\sum_{\mathbf{x}} \mathbf{q}(\mathbf{x}|\mathbf{u}) d_1(u_i, x_i) \leq D_1, \quad \forall \mathbf{u} \in \mathcal{A}^n, 1 \leq i \leq n \quad (19)$$

The single-letter component Q_* of \mathbf{Q}_* is given by the minimizer of $F(\rho, Q, W)$ subject to the same constraints applied to the case $n = 1$, i.e.,

$$\min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} F(\rho, \mathbf{Q}, W^n) = \left[\min_{Q \in \mathcal{Q}(D_1)} F(\rho, Q, W) \right]^n. \quad (20)$$

The proof of Theorem 3 is sketched in Section 6.

Comment: The fact that a memoryless input maximizes the random coding error exponent of a memoryless channel, in the absence of distortion constraints, is well-known as the *parallel channels theorem* due to Gallager [5, Theorem 5],[6, p. 149, Example 4]. Theorem 3 tells us that this continues to be the case in the presence of distortion constraints.

We are now ready to characterize the minimax–maximin game between the information hider and the attacker.

Theorem 4 For a memoryless source $\mathbf{P} = P^n$,

$$\begin{aligned} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}) &= \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}) \\ &= 2^{-nE(R, D_1, D_2)}, \end{aligned} \quad (21)$$

where

$$\begin{aligned} E(R, D_1, D_2) &= \max_{0 \leq \rho \leq 1} \max_{Q \in \mathcal{Q}(D_1)} \min_{W \in \mathcal{W}(D_2)} E_r(\rho, Q, W) \\ &= \max_{0 \leq \rho \leq 1} \min_{W \in \mathcal{W}(D_2)} \max_{Q \in \mathcal{Q}(D_1)} E_r(\rho, Q, W). \end{aligned} \quad (22)$$

Moreover, both the minimax and the maximin of eq. (21) are attained by a saddle point $(\mathbf{Q}_*, \mathbf{W}_*) = (Q_*^n, W_*^n)$, where (Q_*, W_*) is a saddle point of $\min_{0 \leq \rho \leq 1} B(\rho, Q, W)$, or equivalently, a saddle point (Q_ρ, W_ρ) of $B(\rho, Q, W)$ for the value of ρ that minimizes $B(\rho, Q_\rho, W_\rho)$.

Proof. For every fixed ρ , $F(\rho, \mathbf{Q}, \mathbf{W})$, and hence also $B(\rho, \mathbf{Q}, \mathbf{W})$, is convex in \mathbf{Q} [5, Theorem 4] and concave in \mathbf{W} (Section 5, Lemma 1). Also, $\mathcal{Q}_n(D_1)$ and $\mathcal{W}_n(D_2)$ are both compact convex sets. Therefore,

$$\max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} B(\rho, \mathbf{Q}, \mathbf{W}) = \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} B(\rho, \mathbf{Q}, \mathbf{W}) \quad (23)$$

and hence there exists a saddle point [8, p. 229, Theorem 6.2.9(ii)] $(\mathbf{Q}_\rho, \mathbf{W}_\rho)$ depending on ρ , i.e.,

$$B(\rho, \mathbf{Q}, \mathbf{W}) \leq B(\rho, \mathbf{Q}_\rho, \mathbf{W}_\rho) \leq B(\rho, \mathbf{Q}, \mathbf{W}_\rho) \quad (24)$$

for every $(\mathbf{Q}, \mathbf{W}) \in \mathcal{Q}_n(D_1) \times \mathcal{W}_n(D_2)$. Taking the minimum over ρ , we get

$$\min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}) \leq \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}_\rho, \mathbf{W}_\rho) \leq \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}_\rho). \quad (25)$$

Since this is true for every $(\mathbf{Q}, \mathbf{W}) \in \mathcal{Q}_n(D_1) \times \mathcal{W}_n(D_2)$, one may maximize the l.h.s. over $\mathbf{W} \in \mathcal{W}_n(D_2)$ and minimize the r.h.s. over $\mathbf{Q} \in \mathcal{Q}_n(D_1)$, which leads to

$$\max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}) \leq \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}_\rho, \mathbf{W}_\rho) \leq \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}_\rho). \quad (26)$$

The left-most side of eq. (26) is lower bounded by

$$\min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W})$$

whereas the right-most side of eq. (26) is upper bounded by

$$\max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}).$$

On the other hand, since minmax is never smaller than maxmin, all the inequalities must be equalities, i.e.,

$$\begin{aligned} \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}) &= \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}_\rho, \mathbf{W}_\rho) \\ &= \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} \min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W}), \end{aligned} \quad (27)$$

establishing the first equality in eq. (21). Moreover, eq. (27) tells us that the saddle point $(\mathbf{Q}_*, \mathbf{W}_*)$ is given by $(\mathbf{Q}_{\rho_0}, \mathbf{W}_{\rho_0})$, where ρ_0 minimizes $B(\rho, \mathbf{Q}_\rho, \mathbf{W}_\rho)$. Let (Q_ρ, W_ρ) be the saddle point of $B(\rho, Q, W)$ corresponding to $n = 1$. By Theorem 1, for every $\mathbf{W} \in \mathcal{W}_n(D_2)$,

$$\begin{aligned} B(\rho, \mathbf{Q}_\rho^n, \mathbf{W}) &\leq \max_{\mathbf{W} \in \mathcal{W}_n(D_2)} B(\rho, \mathbf{Q}_\rho^n, \mathbf{W}) \\ &= \left[\max_{W \in \mathcal{W}(D_2)} B(\rho, Q_\rho, W) \right]^n \\ &= B(\rho, Q_\rho, W_\rho)^n \\ &= B(\rho, Q_\rho^n, W_\rho^n). \end{aligned} \quad (28)$$

On the other hand, by Theorem 3, for every $\mathbf{Q} \in \mathcal{Q}_n(D_1)$,

$$\begin{aligned} B(\rho, \mathbf{Q}, \mathbf{W}_\rho^n) &\geq \min_{\mathbf{Q} \in \mathcal{Q}_n(D_1)} B(\rho, \mathbf{Q}, \mathbf{W}_\rho^n) \\ &= \left[\min_{Q \in \mathcal{Q}(D_1)} B(\rho, Q, W) \right]^n \\ &= B(\rho, Q_\rho, W_\rho)^n \\ &= B(\rho, Q_\rho^n, W_\rho^n). \end{aligned} \quad (29)$$

Thus, (Q_ρ^n, W_ρ^n) is a saddle point of $B(\rho, \mathbf{Q}, \mathbf{W})$ for every ρ . In particular, for $\rho = \rho_0$, eq. (27) tells us that $(Q_{\rho_0}^n, W_{\rho_0}^n)$ achieves both the minimax and the maximin of $\min_{0 \leq \rho \leq 1} B(\rho, \mathbf{Q}, \mathbf{W})$. The

two expressions of $E(R, D_1, D_2)$ (the first one has been shown already in Theorem 1) follow now immediately from eq. (27) applied to the case $n = 1$. \square

5 Proof of Theorem 1

The proof involves three auxiliary results stated below as lemmas. The first important property of F is concavity w.r.t. \mathbf{W} .

Lemma 1 *The functional $F(\rho, \mathbf{Q}, \mathbf{W})$ is concave in \mathbf{W} for fixed \mathbf{Q} and ρ .*

Proof. It would be sufficient to show that for every $(u, y) \in \mathcal{A} \times \mathcal{A}$, each summand of the r.h.s. of eq. (5) is concave in \mathbf{W} . Specifically, let \mathbf{W}_1 and \mathbf{W}_2 be two arbitrary channels and let α be an arbitrary number in $(0, 1)$. We wish to show then that

$$\begin{aligned} & \left\{ \sum_{\mathbf{x}} q(\mathbf{x}|\mathbf{u}) [\alpha w_1(\mathbf{y}|\mathbf{x}) + (1 - \alpha)w_2(\mathbf{y}|\mathbf{x})]^{1/(1+\rho)} \right\}^{1+\rho} \\ & \geq \alpha \left[\sum_{\mathbf{x}} q(\mathbf{x}|\mathbf{u}) w_1(\mathbf{y}|\mathbf{x})^{1/(1+\rho)} \right]^{1+\rho} + (1 - \alpha) \left[\sum_{\mathbf{x}} q(\mathbf{x}|\mathbf{u}) w_2(\mathbf{y}|\mathbf{x})^{1/(1+\rho)} \right]^{1+\rho} \\ & = \left\{ \sum_{\mathbf{x}} [\alpha q(\mathbf{x}|\mathbf{u})^{1+\rho} w_1(\mathbf{y}|\mathbf{x})]^{1/(1+\rho)} \right\}^{1+\rho} + \left\{ \sum_{\mathbf{x}} [(1 - \alpha)q(\mathbf{x}|\mathbf{u})^{1+\rho} w_2(\mathbf{y}|\mathbf{x})]^{1/(1+\rho)} \right\}^{1+\rho}. \end{aligned} \quad (30)$$

Now, Minkowski's inequality (see, e.g., [6, p. 523, Problem 4.15(g)]) states that for a set of non-negative numbers $\{a_{jk}\}$, $1 \leq j \leq J$, $1 \leq k \leq K$, and $r > 1$,

$$\left[\sum_j \left(\sum_k a_{jk} \right)^{1/r} \right]^r \geq \sum_k \left(\sum_j a_{jk}^{1/r} \right)^r. \quad (31)$$

It is easy then to see that eq. (30) is obtained as a special case of Minkowski's inequality with $r = 1 + \rho$, $K = 2$, $j = \mathbf{x}$, $a_{j1} = \alpha q(\mathbf{x}|\mathbf{u})^{1+\rho} w_1(\mathbf{y}|\mathbf{x})$, and $a_{j2} = (1 - \alpha)q(\mathbf{x}|\mathbf{u})^{1+\rho} w_2(\mathbf{y}|\mathbf{x})$. \square

Consider now an auxiliary maximization problem, henceforth referred to as *Problem A*, which is defined as follows. Let $\mathcal{N} = \{1, 2, \dots, n\}$ and let S be a subset of \mathcal{N} , where the members are

listed in increasing order, e.g., $\{1, 2, 5\}$, $\{3, 4, 100, 200\}$, etc. Let $n(S)$ denote the cardinality of S , and let \mathcal{S} denote the collection of all $2^n - 1$ nonempty subsets of \mathcal{N} . Finally, let $\theta > 0$ be a given constant. Then, Problem A is defined as follows.

Problem A:

$$\max_{\mathbf{W}} F(\rho, \mathbf{Q}, \mathbf{W}) \quad (32)$$

subject to

$$\mathbf{w}(\mathbf{y}|\mathbf{x}) \geq 0 \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{A}^n \quad (33)$$

$$\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) = 1 \quad \forall \mathbf{x} \in \mathcal{A}^n \quad (34)$$

$$\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) \exp\left[\theta \sum_{i \in S} d_2(x_i, y_i)\right] \leq \exp[\theta n(S) D_2], \quad \forall \mathbf{x} \in \mathcal{A}^n, S \in \mathcal{S}. \quad (35)$$

The reason for considering Problem A is the fact that it is easily shown to be solved by a product form channel. This will follow from a simple inspection of the Kuhn–Tucker conditions (similarly as in [5], [6, Theorems 4.5.1, 5.6.5]). The original problem of interest will be approached by letting $\theta \rightarrow 0$.

Lemma 2 *For a given P , Q , ρ , and D_2 , let $W_* = \{w_*(y|x)\}$ maximize $F(\rho, Q, W)$ subject to the constraints*

$$w(y|x) \geq 0 \quad \forall x, y \in \mathcal{A} \quad (36)$$

$$\sum_y w(y|x) = 1 \quad \forall x \in \mathcal{A} \quad (37)$$

$$\sum_y w(y|x) e^{\theta d_2(x,y)} \leq e^{\theta D_2} \quad \forall x \in \mathcal{A}. \quad (38)$$

Then, for $\mathbf{P} = P^n$ and $\mathbf{Q} = Q^n$, the channel $\mathbf{W} = W_^n$ solves Problem A.*

Proof. First observe that the solution to Problem A will be unaffected if the equality constraint (34) will be replaced by the inequality constraint $\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) \leq 1$ for all \mathbf{x} . This is obvious because for any \mathbf{W} with $\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) < 1$ for some \mathbf{x} , the components $\mathbf{w}(\mathbf{y}|\mathbf{x})$ can be scaled by a factor larger than one without violating the constraint, thereby increasing F . Therefore, the maximum of F must be attained for $\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) = 1$ for all \mathbf{x} . Let us refer to this modification³ of Problem A as Problem A'. Next observe that Problem A' is a convex program (see, e.g., [3, Sect. 4.5]) with inequality constraints only. Therefore, the Kuhn–Tucker conditions are both necessary and sufficient for optimality [3, Theorems 4.38, 4.39]. Specifically, let

$$\begin{aligned} V(\mathbf{Q}, \mathbf{W}, \mathbf{x}, \mathbf{y}) &\triangleq \frac{\partial F(\rho, \mathbf{Q}, \mathbf{W})}{\partial \mathbf{w}(\mathbf{y}|\mathbf{x})} \\ &= \mathbf{w}(\mathbf{y}|\mathbf{x})^{-\rho/(1+\rho)} \sum_{\mathbf{u}} \mathbf{p}(\mathbf{u})\mathbf{q}(\mathbf{x}|\mathbf{u}) \left[\sum_{\mathbf{x}'} \mathbf{q}(\mathbf{x}'|\mathbf{u})\mathbf{w}(\mathbf{y}|\mathbf{x}')^{1/(1+\rho)} \right]^{\rho}. \end{aligned} \quad (39)$$

The Kuhn–Tucker conditions for a channel \mathbf{W} to solve Problem A' are that there exist non-negative constants $\{\boldsymbol{\mu}(\mathbf{x})\}$ and $\{\boldsymbol{\lambda}(\mathbf{x}, S)\}$, $\mathbf{x} \in \mathcal{A}^n$, $S \in \mathcal{S}$, such that for every \mathbf{x} and \mathbf{y} ,

$$V(\mathbf{Q}, \mathbf{W}, \mathbf{x}, \mathbf{y}) = \boldsymbol{\mu}(\mathbf{x}) + \sum_{S \in \mathcal{S}} \boldsymbol{\lambda}(\mathbf{x}, S) \exp[\theta \sum_{i \in S} d_2(x_i, y_i)], \quad \mathbf{w}(\mathbf{y}|\mathbf{x}) > 0 \quad (40)$$

$$V(\mathbf{Q}, \mathbf{W}, \mathbf{x}, \mathbf{y}) \leq \boldsymbol{\mu}(\mathbf{x}) + \sum_{S \in \mathcal{S}} \boldsymbol{\lambda}(\mathbf{x}, S) \exp[\theta \sum_{i \in S} d_2(x_i, y_i)], \quad \mathbf{w}(\mathbf{y}|\mathbf{x}) = 0. \quad (41)$$

Let $W_* = \{w_*(y|x), x, y \in \mathcal{A}\}$ be defined as in Lemma 2. Then, by the necessity of the Kuhn–Tucker conditions for $n = 1$, there must exist non-negative constants $\lambda(x)$, and $\mu(x)$, $x \in \mathcal{A}$, such that

$$V(Q, W_*, x_i, y_i) = \mu(x_i) + \lambda(x_i) \exp[\theta d_2(x_i, y_i)], \quad w_*(y_i|x_i) > 0 \quad (42)$$

$$V(Q, W_*, x_i, y_i) \leq \mu(x_i) + \lambda(x_i) \exp[\theta d_2(x_i, y_i)], \quad w_*(y_i|x_i) = 0. \quad (43)$$

³This modification guarantees that the Lagrange multiplier corresponding to the equality constraint is non-negative.

Taking the product of $\{V(Q, W_*, x_i, y_i)\}$ and $\{\mu(x_i) + \lambda(x_i) \exp[\theta d_2(x_i, y_i)]\}$ over $1 \leq i \leq n$, the former gives $V(Q, W_*^n, \mathbf{x}, \mathbf{y})$, whereas and the latter is

$$\prod_i [\mu(x_i) + \lambda(x_i) \exp[\theta d_2(x_i, y_i)]] \triangleq \boldsymbol{\mu}(\mathbf{x}) + \sum_{S \in \mathcal{S}} \boldsymbol{\lambda}(\mathbf{x}, S) \exp[\theta \sum_{i \in S} d_2(x_i, y_i)] \quad (44)$$

with

$$\boldsymbol{\mu}(\mathbf{x}) = \prod_{i=1}^n \mu(x_i) \quad (45)$$

and

$$\boldsymbol{\lambda}(\mathbf{x}, S) = \prod_{i \in S} \lambda(x_i) \prod_{i \in S^c} \mu(x_i). \quad (46)$$

Clearly, since $\mu_i(x)$ and $\lambda_i(x)$ are all non-negative, so are $\boldsymbol{\mu}(\mathbf{x})$ and $\boldsymbol{\lambda}(\mathbf{x}, S)$. In addition, if $w_*(y_i|x_i) > 0$ for all i (and hence also $\mathbf{w}_*(\mathbf{y}|\mathbf{x}) = \prod_i w_*(y_i|x_i) > 0$), then

$$V(Q, \mathbf{W}, \mathbf{x}, \mathbf{y}) = \boldsymbol{\mu}(\mathbf{x}) + \sum_S \boldsymbol{\lambda}(\mathbf{x}, S) \exp[\theta \sum_{i \in S} d_2(x_i, y_i)],$$

otherwise

$$V(Q, \mathbf{W}, \mathbf{x}, \mathbf{y}) \leq \boldsymbol{\mu}(\mathbf{x}) + \sum_S \boldsymbol{\lambda}(\mathbf{x}, S) \exp[\theta \sum_{i \in S} d_2(x_i, y_i)].$$

The conclusion is, therefore, that the Lagrange multiplier values given in eqs. (45) and (46) satisfy the Kuhn–Tucker conditions together with the vector channel $\mathbf{W} = W_*^n$. Thus, by the sufficiency of the Kuhn–Tucker conditions, this product form channel solves Problem A', or equivalently, Problem A. \square .

Let us now define Problem B similarly as Problem A, but with the the distortion constraints (35) replaced by

$$\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) d_2(x_i, y_i) \leq D_2, \quad \forall \mathbf{x} \in \mathcal{A}^n, 1 \leq i \leq n. \quad (47)$$

In other words, Problem B is the original problem under discussion.

Lemma 3 Let $\mathbf{P} = P^n$, \mathbf{Q} , ρ , and D_2 be given, and define the function $F_\theta(D_2)$ as the maximum value of $F(\rho, Q^n, \mathbf{W})$ subject to the constraints of Problem A. Let $F_0(D_2)$ be defined as the maximum value of $F(\rho, Q^n, \mathbf{W})$ subject to the constraints of Problem B. Then,

(a) The function $F_0(D_2)$ is concave.

(b) $\lim_{\theta \rightarrow 0} F_\theta(D_2) = F_0(D_2)$.

(c) The value $F_0(D_2)$ is attained by a memoryless channel W_*^n , where W_* solves Problem B for $n = 1$ w.r.t. the distortion constraints $\sum_y w(y|x)d_2(x,y) \leq D_2$ for all $x \in \mathcal{A}$.

Proof. Beginning from part (a), let D'_2 and D''_2 denote two distortion levels and for $\alpha \in (0, 1)$, let $D_2 = \alpha D'_2 + (1 - \alpha)D''_2$. Let \mathbf{W}'_* , \mathbf{W}''_* , and \mathbf{W}_* solve Problem B for D'_2 , D''_2 , and D_2 , respectively.

Now, we have

$$\begin{aligned}
\alpha F_0(D'_2) + (1 - \alpha)F_0(D''_2) &= \alpha F(\rho, \mathbf{Q}, \mathbf{W}'_*) + (1 - \alpha)F(\rho, \mathbf{Q}, \mathbf{W}''_*) \\
&\leq F(\rho, \mathbf{Q}, \alpha \mathbf{W}'_* + (1 - \alpha)\mathbf{W}''_*) \\
&\leq F(\rho, \mathbf{Q}, \mathbf{W}_*) \\
&= F_0(D_2),
\end{aligned} \tag{48}$$

where the first inequality follows from the concavity of F (Lemma 1), and the second inequality results from the fact that $\alpha \mathbf{W}'_* + (1 - \alpha)\mathbf{W}''_*$ satisfies the constraints of Problem B with distortion matrix D_2 .

To prove (b), observe first that the set of individual coordinate distortion constraints given in eq. (47) is equivalent to the set of constraints

$$\sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x}) \left[\sum_{i \in S} d_2(x_i, y_i) \right] \leq n(S)D_2, \quad \forall \mathbf{x} \in \mathcal{A}^n, S \in \mathcal{S}. \tag{49}$$

This is true because on the one hand, the subsets $\{1\}, \{2\}, \dots, \{n\}$ are all in \mathcal{S} , and on the other hand, given the individual coordinate distortion constraints, the remaining distortion constraints are all redundant. Now, for every $(d_1, \dots, d_n) \in [0, D_{\max}]^n$ and $S \in \mathcal{S}$, $e^{\theta \sum_{i \in S} d_i} \geq 1 + \theta \sum_{i \in S} d_i$, and for $0 < \theta \leq 1$, we also have

$$\begin{aligned}
e^{\theta \sum_{i \in S} d_i} &\leq 1 + \theta \sum_{i \in S} d_i + \sum_{k \geq 2} \frac{(\theta n D_{\max})^k}{k!} \\
&= 1 + \theta \sum_{i \in S} d_i + \theta^2 \sum_{k \geq 2} \frac{\theta^{k-2} (n D_{\max})^k}{k!} \\
&\leq 1 + \theta \sum_{i \in S} d_i + \theta^2 \sum_{k \geq 2} \frac{(n D_{\max})^k}{k!} \\
&= 1 + \theta \sum_{i \in S} d_i + \theta^2 (e^{n D_{\max}} - n D_{\max} - 1) \\
&= 1 + \theta \sum_{i \in S} d_i + C \theta^2,
\end{aligned} \tag{50}$$

where $C = e^{n D_{\max}} - n D_{\max} - 1$ is a constant (for a given n) that depends on neither θ nor $\{d_i\}$. Let $\mathcal{R}_\theta(D_2)$ denote the set of channels \mathbf{W} that satisfy the constraints of Problem A and let $\mathcal{R}_0(D_2)$ denote the set of channels that satisfy the constraints of Problem B. By using the above upper and lower bounds on the exponential function, it is readily seen that for $\theta \in (0, 1]$,

$$\mathcal{R}_0(D_2 - \theta C) \subseteq \mathcal{R}_\theta(D_2) \subseteq \mathcal{R}_0(D_2 + \theta C), \tag{51}$$

and therefore,

$$F_0(D_2 - \theta C) \leq F_\theta(D_2) \leq F_0(D_2 + \theta C). \tag{52}$$

From the concavity of F_0 , proven already in part (a) above, it also follows that F_0 is continuous at any point in $(0, D_{\max})$. Thus, for such D_2 , $\lim_{\theta \rightarrow 0} F_0(D_2 - \theta C) = \lim_{\theta \rightarrow 0} F_0(D_2 + \theta C) = F_0(D_2)$, and assertion (b) follows from eq. (52) by a sandwich argument.

As for part (c), let $\{\theta_k\}$ be a positive sequence that tends to zero, and consider the corresponding sequence of channels $\{\mathbf{W}_k\}$ that achieve $F_{\theta_k}(D_2)$, $k = 1, 2, \dots$. According to Lemma 2, $\mathbf{W}_k = W_k^n$,

where W_k solves the single-letter version of Problem A (i.e., $\max_W F(\rho, Q, W)$ subject to the constraints of Lemma 2). Since the set of finite-alphabet channels is compact, there is a convergent subsequence $\{W_{k_j}\}_{j \geq 1}$. Let W_* denote the limit of that subsequence. Then, from the continuity of F w.r.t. \mathbf{W} ,

$$\begin{aligned} F(\rho, \mathbf{Q}, W_*^n) &= \lim_{j \rightarrow \infty} F(\rho, \mathbf{Q}, W_{k_j}^n) \\ &= \lim_{\theta \rightarrow 0} F_\theta(D_2) \\ &= F_0(D_2), \end{aligned} \tag{53}$$

where the last step follows from part (b). Thus, Problem B is solved by W_*^n . Finally, the fact that W_* solves the single-letter version of Problem B, can be readily seen again from eq. (53), degenerated to the one-dimensional case. \square

6 Proof of Theorem 3

The proof is completely analog to that of Theorem 1. Therefore, we give here a brief sketch only. The fact that $F(\rho, \mathbf{Q}, \mathbf{W})$ is convex w.r.t. \mathbf{Q} is well-known [5, Theorem 4]. Next, define Problem A as

$$\min_{\mathbf{Q}} F(\rho, \mathbf{Q}, W^n) \tag{54}$$

subject to

$$\mathbf{q}(\mathbf{x}|\mathbf{u}) \geq 0 \quad \forall \mathbf{u}, \mathbf{x} \in \mathcal{A}^n \tag{55}$$

$$\sum_{\mathbf{x}} \mathbf{q}(\mathbf{x}|\mathbf{u}) = 1 \quad \forall \mathbf{u} \in \mathcal{A}^n \tag{56}$$

$$\sum_{\mathbf{y}} \mathbf{q}(\mathbf{x}|\mathbf{u}) \exp[-\theta \sum_{i \in S} d_1(u_i, x_i)] \geq \exp[-\theta n(S)D_1], \quad \forall \mathbf{u} \in \mathcal{A}^n, S \in \mathcal{S}. \tag{57}$$

The first step is to show that Problem A is solved by a product form channel $\mathbf{Q} = Q^n$, where Q solves the one-dimensional version of Problem A corresponding to each coordinate. To see this, Problem A is first modified (without affecting the result) to Problem A', defined the same way as Problem A but with the equality constraint replaced by $\sum_{\mathbf{x}} \mathbf{q}(\mathbf{x}|\mathbf{u}) \geq 1$ for all \mathbf{u} . This is again a convex program for which the Kuhn–Tucker (necessary and sufficient) conditions are the existence of non-negative constants $\{\boldsymbol{\mu}(\mathbf{u})\}$, $\{\boldsymbol{\lambda}(\mathbf{u}, S)\}$, $\mathbf{u} \in \mathcal{A}^n$, $S \in \mathcal{S}$, such that

$$U(\mathbf{Q}, \mathbf{W}, \mathbf{u}, \mathbf{x}) = \boldsymbol{\mu}(\mathbf{u}) + \sum_{S \in \mathcal{S}} \boldsymbol{\lambda}(\mathbf{u}, S) \exp[-\theta \sum_{i \in S} d_1(u_i, x_i)], \quad \mathbf{q}(\mathbf{x}|\mathbf{u}) > 0 \quad (58)$$

$$U(\mathbf{Q}, \mathbf{W}, \mathbf{u}, \mathbf{x}) \geq \boldsymbol{\mu}(\mathbf{u}) + \sum_{S \in \mathcal{S}} \boldsymbol{\lambda}(\mathbf{u}, S) \exp[-\theta \sum_{i \in S} d_1(u_i, x_i)], \quad \mathbf{q}(\mathbf{x}|\mathbf{u}) = 0 \quad (59)$$

where

$$U(\mathbf{Q}, \mathbf{W}, \mathbf{u}, \mathbf{x}) = \mathbf{p}(\mathbf{u}) \sum_{\mathbf{y}} \mathbf{w}(\mathbf{y}|\mathbf{x})^{1/(1+\rho)} \left[\sum_{\mathbf{x}'} \mathbf{q}(\mathbf{x}'|\mathbf{u}) \mathbf{w}(\mathbf{y}|\mathbf{x}')^{1/(1+\rho)} \right]^\rho \quad (60)$$

is the partial derivative of F w.r.t. $\mathbf{q}(\mathbf{x}|\mathbf{u})$. Similarly as in the proof of Lemma 2, these conditions give rise to the above-defined product form channel.

Next, define Problem B in the same way as Problem A, but with the set of distortion constraints replaced by

$$\sum_{\mathbf{x}} \mathbf{q}(\mathbf{x}|\mathbf{u}) d_1(u_i, x_i) \leq D_1, \quad \forall \mathbf{u} \in \mathcal{A}^n, 1 \leq i \leq n. \quad (61)$$

Finally, an analogue to Lemma 3 (with the concavity property in part (a) replaced by a convexity property) is then proved in the same way, where the bounds on the exponential function to be used are $e^{-\theta \sum_{i \in S} d_i} \geq 1 - \theta \sum_{i \in S} d_i$ and $e^{-\theta \sum_{i \in S} d_i} \leq 1 - \theta \sum_{i \in S} d_i + \frac{1}{2} \theta^2 (n D_{\max})^2$.

Acknowledgement

Useful discussions with Ronny Roth are greatly appreciated.

References

- [1] R. Anderson, “Stretching the limits of steganography,” in *Information Hiding*, Springer Lecture Notes in Computer Science vol. 1174, pp. 39–48, 1996. Also, available at <http://www.bluespike.com/research/research.html>.
- [2] R. J. Anderson and A. Petitcolas, “On the limits of steganography,” *IEEE J. Select. Areas in Communications*, vol. 16, no. 4, pp. 463–473, May 1998.
- [3] M. Avriel, *Nonlinear Programming Analysis and Methods*. Prentice-Hall, 1976.
- [4] S. Craver, N. Memon, B.-L. Yeo, and M. M. Yeung, “Resolving rightful ownerships with invisible watermarking techniques: limitations, attacks, and implications,” *IEEE J. Select. Areas in Commun.*, vol. 16, no. 4, pp. 573–586, May 1998.
- [5] R. G. Gallager, “A simple derivation of the coding theorem and some applications,” *IEEE Trans. Inform. Theory*, vol. IT–11, no. 1, pp. 3–18, 1965.
- [6] R. G. Gallager, *Information Theory and Reliable Communication*, J. Wiley & Sons, 1968.
- [7] J. A. O’Sullivan, P. Moulin, and J. M. Ettinger, “Information theoretic analysis of Steganography,” *Proc. ISIT ‘98*, p. 297, 1998.
- [8] J. Stoer and C. Witzgall, *Convexity and Optimization in Finite Dimensions I*, Springer-Verlag, New York, 1970.