One shot approach to lossy source coding under average distortion constraints

Nir Elkayam Meir Feder Department of Electrical Engineering - Systems Tel-Aviv University, Israel Email: nirelkayam@post.tau.ac.il, meir@eng.tau.ac.il

Abstract

This paper presents a one shot analysis of the lossy compression problem under average distortion constraints. We calculate the exact expected distortion of a random code. The result is given as an integral formula using a newly defined functional $\tilde{D}(z, Q_Y)$ where Q_Y is the random coding distribution and $z \in [0, 1]$. When we plug in the code distribution as Q_Y , this functional produces the average distortion of the code, thus provide a converse result utilizing the same functional. Two alternative formulas are provided for $\tilde{D}(z, Q_Y)$, the first involves a supremum over some auxiliary distribution Q_X which has resemblance to the channel coding meta-converse and the other involves an infimum over channels which resemble the well known Shannon distortion-rate function.

I. INTRODUCTION

The single shot approach aims to find informational quantities that govern the optimal performance of an operational problems of interest, *e.g.*, channel coding and lossy compression. In both cases, the problem settings pose a random object that we want to control. In the channel coding problem this is the random channel which abstracts the medium we want to use to enable a reliable communication. In the lossy compression this is the source we want to "compress" to a minimum number of bits subject to a distortion constraint. The single shot approach tries to solve the problem by providing achievable and converse bounds without any assumption on the random object.

A "good" solution should have the following properties:

- 1) **Tightness**: The relation between the achievable and converse bounds should be clarifies and quantified. Preferably, the gap between the bounds should be "small".
- 2) Computation: The bounds should be computable. Since we generally deal with a high dimensional problem space for which the exact description might not even be feasible, we relax the computability to convexity, *i.e.*, the bounds should be presented as a minimization of some convex function on some convex domain. For such a problems, symmetries might solve the problem entirely or substantially reduce the effective size, see *e.g.*[1, Theorem 20].
- 3) Generalization: The bounds can be relaxed to other known bounds.

In this paper we deal with the lossy source coding problem. In [2] we presented a general approach to the one-shot coding problem. In this paper we borrow and extends ideas from [2] and provide a novel analysis of the lossy compression problem. We derive an achievable bound using random coding and a corresponding converse bound. Both bounds are given in term of a newly defined functional $\tilde{D}(z, Q_Y)$ where $z \in [0, 1]$ and Q_Y denote a distribution over the reproduction space \mathcal{Y} . The functional $\tilde{D}(z, Q_Y)$ is shown to be convex and has similarity to the channel coding meta-converse [3, Theorem 27].

A. Notation

Throughout this paper, scalar random variables are denoted by capital letters (e.g. X), sample values are denoted by lower case letters (e.g. x) and their alphabets are denoted by their respective calligraphic letters, (e.g. \mathcal{X}).

The set of all distributions (probability mass functions) supported on alphabet \mathcal{Y} is denoted as $\mathscr{P}(\mathcal{Y})$. The set of all conditional distributions (i.e., channels) with the input alphabet \mathcal{X} and the output alphabet \mathcal{Y} is denoted by $\mathscr{P}(\mathcal{Y}|\mathcal{X})$. If X has distribution P_X , we write this as $X \sim P_X$. The uniform probability distribution over [0, 1] is denoted throughout by \mathcal{U} . The probability (expectation) of an event (random variable) \mathcal{A} under the distribution P_X is denoted by $\mathbb{P}_{P_X} \{\mathcal{A}\}$ ($\mathbb{E}_{P_X}(\mathcal{A})$) respectively, e.g. $\mathbb{P}_{P_X} \{X \ge \alpha\}$ and $\mathbb{E}_{P_X} (f(X))$. In some cases, we abbreviate the notation and write $P_X \{\mathcal{A}\}$ instead of $\mathbb{P}_{P_X} \{\mathcal{A}\}$, e.g. $P_X (X \ge \alpha) = \mathbb{P}_{P_X} \{X \ge \alpha\}$. In some cases, we write $\mathbb{E}_{\mu} (f(X))$ where μ is σ -finite measure and not a probability measure.

II. PROBLEM SETTING

Let X denote the random variable on \mathcal{X} , representing the source we want to compress. The elements of \mathcal{X} are the input symbols. Denote by P_X the distribution of X. Let \mathcal{Y} denote the set of reproduction symbols. Let $d : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}^+$ denote the distortion function. Any subset $\mathcal{C} \subset \mathcal{Y}$ is a **code** and the *average distortion* associated with this code is:

$$D(P_X, \mathcal{C}) \triangleq \mathbb{E}_{P_X} \left(\min_{y \in \mathcal{C}} \left\{ d(X, y) \right\} \right)$$
(1)

Let:

$$D(P_X, R) \triangleq \min_{\mathcal{C} \subset \mathcal{V} \cap \mathcal{C}^+ = R} D(P_X, \mathcal{C})$$
(2)

denote the optimal distortion-rate function. The goal is to find upper and lower bounds on $D(P_X, R)$.

Throughout this paper we assume that both \mathcal{X} and \mathcal{Y} are finite sets. Thus the distribution P_X is discrete and the distortion is bounded by some d_{max} . The results in this paper can be extended quite straight forwardly to rather general alphabets \mathcal{X} , \mathcal{Y} and appropriate σ -algebras, as long as the probability distribution P_X is well defined. The boundedness of the distortion can be relaxed to the following: There exist a "small" finite set $\mathcal{R} \subset \mathcal{Y}$ such that $\mathbb{E}_{P_X}(\min_{u \in \mathcal{U}} d(X, y)) = d_{max} < \infty$.

III. ACHIEVABILITY BOUND

For the achievable argument we use the random coding approach. Let $Q_Y \in \mathscr{P}(\mathcal{Y})$ denote a given distribution on \mathcal{Y} . A random code of rate R with $M = e^R + 1$ codewords is $C = \{Y_0, \ldots, Y_{e^R}\}$ where each Y_i is drawn from Q_Y independently of the other code words. The average distortion of the random code is:

$$D(P_X, Q_Y, R) \triangleq \mathbb{E}_C \left(D(P_X, C) \right) \tag{3}$$

A. pairwise correct probability

In [2] the pairwise error probability was defined as the random variable representing the probability of error given the sent and received symbols. Here, we define the *pairwise correct probability*, which represent the probability of drawing a reproduction symbol that is better than a given reproduction symbol.

Definition 1. For $x \in \mathcal{X}$, $y \in \mathcal{Y}$ and $u \in [0, 1]$ Let:

$$p_{c,x,y,u} \triangleq Q_Y \left\{ d(x,Y) < d(x,y) \right\} + u \cdot Q_Y \left\{ d(x,Y) = d(x,y) \right\}$$

$$(4)$$

The pairwise correct decoding probability is the random variable: $p_{c,x,y,U}$ where $U \sim \mathcal{U}$ is uniform over [0, 1].

The following proposition summaries the properties we need about the pairwise correct decoding $p_{c,x,y,u}$:

Proposition 1.

1) For any $w \in [0, 1]$ and x there exist y and τ such that:

$$v = p_{c,x,y,\tau}.$$

- 2) $d(x, y_1) < d(x, y_2) \Rightarrow p_{c,x,y_1,U_1} \le p_{c,x,y_2,U_2}$. If $Q_Y(y_1) > 0$ or $Q_Y(y_2) > 0$ then $p_{c,x,y_1,U_1} < p_{c,x,y_2,U_2}$ with probability 1.
- 3) $p_{c,x,Y,U} \sim W$ where $Y \sim Q_Y$ and U, W are uniform over [0, 1].

Proof. To prove (1) note that there must exist a y such that:

$$Q_Y \{ d(x, Y) < d(x, y) \} \le w \le Q_Y \{ d(x, Y) \le d(x, y) \}$$

If $Q_Y \{d(x, Y) = d(x, y)\} = 0$ then we are done with any τ . If $Q_Y \{d(x, Y) = d(x, y)\} \neq 0$ then

$$\tau = \frac{w - Q_Y \left\{ d(x, Y) < d(x, y) \right\}}{Q_Y \left\{ d(x, Y) = d(x, y) \right\}} \le 1$$

satisfies the requirement. To prove (2):

$$p_{c,x,y_1,U_1} = Q_Y \{ d(x,Y) < d(x,y_1) \} + U_1 \cdot Q_Y \{ d(x,Y) = d(x,y_1) \} \stackrel{(a)}{\leq} Q_Y \{ d(x,Y) \leq d(x,y_1) \} \leq Q_Y \{ d(x,Y) < d(x,y_2) \} \stackrel{(b)}{\leq} Q_Y \{ d(x,Y) < d(x,y_2) \} + U_2 \cdot Q_Y \{ d(x,Y) = d(x,y_2) \} = p_{c,x,y_2,U_2}$$

If $Q_Y(y_1) > 0$ then $Q_Y \{ d(x, Y) = d(x, y_1) \} > 0$ thus we have strict inequality in (a). If $Q_Y(y_2) > 0$ we have strict inequality in (b). To prove (3) let y_w and τ_w be such that:

$$w = Q_Y \{ d(x, Y) < d(x, y_w) \} + \tau_w \cdot Q_Y \{ d(x, Y) = d(x, y_w) \}$$

Then:

$$\begin{aligned} & \mathbb{P}_{Q_Y} \left\{ p_{c,x,Y,U} < w \right\} \\ &= \sum_{y:d(x,y) < d(x,y_w)} Q_Y(y) \mathbb{P}_{Q_Y} \left\{ p_{c,x,y,U} < w \right\} \\ &+ \sum_{y:d(x,y) = d(x,y_w)} Q_Y(y) \mathbb{P}_{Q_Y} \left\{ p_{c,x,y,U} < w \right\} \\ &+ \sum_{y:d(x,y) > d(x,y_w)} Q_Y(y) \mathbb{P}_{Q_Y} \left\{ p_{c,x,y,U} < w \right\} \end{aligned}$$

If $d(x, y) \leq d(x, y_w)$ and $Q_Y(y) > 0$ then:

$$p_{c,x,y,U} \leqslant w = p_{c,x,y_w,T}$$

according to (2). Thus the first sum gives $Q_Y \{ d(x, Y) < d(x, y_w) \}$ and the last sum vanishes. If $d(x, y) = d(x, y_w)$ and $Q_Y(y) > 0$ then:

$$\mathbb{P}_{Q_Y} \left\{ p_{c,x,y,U} < w \right\} = \mathbb{P}_{Q_Y} \left\{ p_{c,x,y_w,U} < p_{c,x,y_w,\tau_w} \right\}$$
$$= \mathbb{P} \left\{ U < \tau_w \right\} = \tau_w$$

hence the middle sum gives $\tau_w \cdot Q_Y \{ d(x, Y) = d(x, y_w) \}$. Combined:

$$\mathbb{P}_{Q_Y} \{ p_{c,x,Y,U} < w \} = Q_Y \{ d(x,Y) < d(x,y_w) \} \\ + \tau_w \cdot Q_Y \{ d(x,Y) = d(x,y_w) \} \\ = w$$

B. Random coding performance

Proposition 1 suggests that for any fixed x we have a correspondence between the elements of \mathcal{Y} and the sub interval

$$[Q_Y \{ d(x, Y) < d(x, y) \}, Q_Y \{ d(x, Y) \le d(x, y) \}] \subset [0, 1]$$

given by $y \iff p_{c,x,y,U}$. We define the distortion function d on $\mathcal{X} \times [0,1]$ according to:

$$d(x,u) \triangleq d(x,y), \quad u = p_{c,x,y,\tau} \tag{5}$$

This correspondence is well defined almost everywhere with respect to the pair (Q_Y, U) since if $d(x, y_1) \neq d(x, y_2)$, and $Q_Y(y_1) > 0$ or $Q_Y(y_2) > 0$ the support of p_{c,x,y_1,U_1} and p_{c,x,y_2,U_2} do not overlap with probability 1. Moreover, this mapping is order preserving, *i.e.*,

$$d(x, y_1) < d(x, y_2) \Rightarrow p_{c, x, y_1, U_1} \le p_{c, x, y_2, U_2}$$
(6)

The following result provides an exact formula for the average distortion of random code.

Theorem 2 (Exact performance of random coding). The average distortion of random code with $M = e^R + 1$ codewords $\{Y_i\}$ drawn from Q_Y is given by:

$$\mathbb{E}_{P_X,\{Y_i\}}\left(\min d(X,Y_i)\right) = \int_0^1 \tilde{D}(w,Q_Y)G'_M(w)dw$$
(7)

where: $\tilde{D}(w, Q_Y) = w^{-1} \cdot \mathbb{E}_{P_X \times Q_Y} \left(d(X, Y) \cdot \mathbb{1}_{\{p_{c,X,Y,U} \le w\}} \right)$ and $G_M(w) = -(1-w)^{M-1}((M-1)w+1)$. Corollary 3. For any $\lambda < R$:

$$\mathbb{E} \left(\min d(X, Y_i) \right) \le \tilde{D} \left(e^{-(R-\lambda)}, Q_Y \right) \\ + \left(\tilde{D} \left(1, Q_Y \right) - \tilde{D} \left(e^{-(R-\lambda)}, Q_Y \right) \right) \cdot e^{-e^{\lambda}} (e^{\lambda} + 1) \\ \le \tilde{D} \left(e^{-(R-\lambda)}, Q_Y \right) + d_{max} \cdot e^{-e^{\lambda}} (e^{\lambda} + 1)$$

For $z \in [\tilde{D}(0, Q_Y), \tilde{D}(1, Q_Y)]$ let:

$$\tilde{D}^{-1}(z, Q_Y) \triangleq \inf \left\{ w \in [0, 1] : \tilde{D}(w, Q_Y) \ge z \right\}$$

and:

$$\tilde{R}(z, Q_Y) \triangleq -\log \tilde{D}^{-1}(z, Q_Y) \tag{8}$$

 $\tilde{R}(z, Q_Y)$ is the "rate distortion" function associated with the prior distribution Q_Y .

Corollary 4. For any Q_Y , let $d_{req} \in (\tilde{D}(0, Q_Y), \tilde{D}(1, Q_Y))$ denote a desired distortion level. There exist a code with distortion level d_{req} and rate R such that:

$$\begin{split} R &\leq \min_{z < d_{req}} \tilde{R}(z, Q_Y) + f^{-1} \left(\frac{d_{req} - z}{\tilde{D}(1, Q_Y) - z} \right) \\ &\leq \min_{z < d_{req}} \tilde{R}(z, Q_Y) + g \left(\frac{\tilde{D}(1, Q_Y) - z}{d_{req} - z} \right) \end{split}$$
where $f(t) = e^{-e^t} (e^t + 1)$ and $g(x) = \log \log(x) + \log \left(\frac{2}{3} \right) + \log \left(1 + \sqrt{1 + 9 \left(2 \log(x) \right)^{-1}} \right)$

Proof of theorem 2:. To calculate the average distortion of a random code with M codewords we proceed as follow:

$$\mathbb{E} \left(\min d(X, Y_i) \right) \stackrel{(a)}{=} \mathbb{E}_{P_X} \left(\mathbb{E} \left(\min d(X, Y_i) | X \right) \right) \\ \stackrel{(b)}{=} \mathbb{E}_{P_X} \left(\mathbb{E} \left(\min \left\{ \tilde{d}(X, p_{c,X,Y_i,U_i}) | X \right\} \right) \right) \\ \stackrel{(c)}{=} \mathbb{E}_{P_X} \left(\mathbb{E} \left(\tilde{d}(X, \min_i \{ p_{c,X,Y_i,U_i} \} \right) \right) | X \right) \\ \stackrel{(d)}{=} \mathbb{E}_{P_X} \left(\mathbb{E}_{W_M} \left(\tilde{d}(X, W_M) | X \right) \right) \\ \stackrel{(e)}{=} \mathbb{E}_{W_M} \left(\mathbb{E}_{P_X} \left(\tilde{d}(X, W_M) | W_M \right) \right) \\ \stackrel{(f)}{=} \int_0^1 \mathbb{E}_{P_X} \left(\tilde{d}(X, w) \right) f_M(w) dw \\ \stackrel{(g)}{=} f_M(1) \tilde{D}_1(1) - f_M(0) \tilde{D}_1(0) \\ - \int_0^1 \tilde{D}_1(w) f'_M(w) dw \\ \stackrel{(h)}{=} M \cdot (M-1) \int_0^1 \tilde{D}_2(w) w (1-w)^{M-2} dw \\ \stackrel{(i)}{=} \int_0^1 \tilde{D}_2(w) G'_M(w) dw$$

where W_M is the minimum of M independent uniform random variables, $\tilde{D}_1(w) = \int_0^w \mathbb{E}_X \left(\tilde{d}(X,z) \right) dz$ and $\tilde{D}_2(w) = w^{-1} \cdot \tilde{D}_1(w)$.

- (a) is the law of total expectation with respect to X.
- (b) follow since $d(x, y) = d(x, p_{c,x,y,U})$ according to (5).
- (c) follow since $p_{c,x,y,U}$ preserve the order induce by \tilde{d} according to (6).
- (d) follow since p_{c,x,Y_i,U_i} are all independent and uniform over [0,1].
- (e) is again the law of total expectation with respect to X and W_M .
- (f) is the expectation according to the p.d.f f_M^{1} . (see (18) in the appendix).
- (g) is integration by parts.
- (h) follow since $D_1(0) = f(1) = 0$ and:

$$f'_M(w) = -M \cdot (M-1)(1-w)^{M-2}$$

• (i) follow since $G'_M(w) = M(M-1)w(1-w)^{M-2}$.

¹The p.d.f of W_M is $f_M(w) = M \cdot (1-u)^{M-1}$

Finlay:

$$\begin{split} \tilde{D}_1(w) &= \int_0^w \mathbb{E}_{P_X} \left(\tilde{d}(X, z) \right) dz \\ &\stackrel{(a)}{=} \mathbb{E}_{P_X, W} \left(\tilde{d}(X, W) \cdot \mathbb{1}_{\{W \le w\}} \right) \\ &\stackrel{(b)}{=} \mathbb{E}_{P_X \times Q_Y} \left(\tilde{d}(X, p_{c, X, Y, U}) \cdot \mathbb{1}_{\{p_{c, X, Y, U} \le w\}} \right) \\ &\stackrel{(c)}{=} \mathbb{E}_{P_X \times Q_Y} \left(d(X, Y) \cdot \mathbb{1}_{\{p_{c, X, Y, U} \le w\}} \right) \\ &= w \cdot \tilde{D}(w, Q_Y) \end{split}$$

where (a) follow since W is uniform over [0, 1]. (b) follow since for each x, $p_{c,x,Y,U}$ is uniform over [0, 1] and (c) is (5). *Proof of corollary 3:* $\tilde{D}(w, Q_Y)$ is an increasing function, Thus:

$$\mathbb{E}\left(\min d(X, Y_i)\right) = \int_0^1 \tilde{D}(w, Q_Y) G'_M(w) dw$$

$$\leq \tilde{D}(u, Q_Y) \int_0^u G'_M(w) dw$$

$$+ \tilde{D}(1, Q_Y) \int_u^1 G'_M(w) dw$$

$$= \tilde{D}(u, Q_Y)$$

$$- \left(\tilde{D}(1, Q_Y) - \tilde{D}(u, Q_Y)\right) G_M(u)$$

since $G_M(0) = -1$, $G_M(1) = 0$ and:

$$-G_M(u) = (1-u)^{M-1}((M-1)u+1)$$

$$\leq e^{-u \cdot e^R}(u \cdot e^R + 1) = e^{-e^{\lambda}}(e^{\lambda} + 1)$$

where $u \cdot e^R = e^{\lambda}$. The second bound follow since $\tilde{D}(1, Q_Y) \leq d_{max}$.

Proof of corollary 4:. Using $d_{req} = \mathbb{E}(\min d(X, Y_i))$ and $z = \tilde{D}(e^{-(R-\lambda)}, Q_Y)$ in corollary 3 we have:

$$d_{req} \le z + (D(1, Q_Y) - z) \cdot f(\lambda)$$

Hence:

$$\lambda \le f^{-1} \left(\frac{d_{req} - z}{\tilde{D}(1, Q_Y) - z} \right)$$

since f is decreasing. The result follow from:

$$\tilde{R}(z, Q_Y) = R - \lambda$$

The second bound follow from the technical lemma 1, given in the appendix.

The standard rate distortion analysis usually employ a "test" channel $W_{Y|X}$ that is used for change of measure during the

achievability proof and also serves as the encoding function when the converse is proved. The following Proposition suggest such a channel to be used in our achievability theorem.

Proposition 5. Let
$$W_{Y|X=x}^w \triangleq w^{-1} \cdot Q_Y \cdot \mathbb{1}_{\{p_{c,x,Y,U} \leq w\}}$$
, i.e.

$$W_{Y|X=x}^{w}(y) = w^{-1} \cdot Q_{Y}(y) \cdot \mathbb{P}_{U} \{ p_{c,x,y,U} \le w \}$$

Then $W_{Y|X=x}^w$ is a probability distribution over \mathcal{Y} and:

$$D(w, Q_Y) = \mathbb{E}_{P_X \times W_{Y|X}^w} \left(d(X, Y) \right)$$

Proof. $W_{Y|X=x}^w$ is a distribution since:

$$w = \mathbb{P}_{Q_Y,U} \{ p_{c,x,Y,U} \le w \}$$
$$= \mathbb{E}_{Q_Y,U} \left(\mathbb{1}_{\{ p_{c,x,Y,U} \le w \}} \right)$$

and:

$$\tilde{D}(w, Q_Y) = \mathbb{E}_{P_X \times Q_Y} \left(d(X, Y) \cdot w^{-1} \cdot \mathbb{1}_{\{p_{c,x,Y,U} \le w\}} \right)$$
$$= \mathbb{E}_{P_X \times W_{Y|X}^w} \left(d(X, Y) \right)$$

IV. CONVERSE BOUND

The channel $W_{Y|X}^w$ that was suggested in proposition 5 with the proper distribution Q_Y and parameter $w = M^{-1}$ can serve as the encoding function as the next proposition show.

Proposition 6. Let $\mathcal{C} \subset \mathcal{Y}$ be a code with M codewords. Let $\tilde{W}_{Y|X}$ denote the optimal encoding of the \mathcal{X} to \mathcal{C} , i.e.

$$\tilde{W}_{Y|X} = \operatorname*{arg\,min}_{c \in \mathcal{C}} \left(d(X, c) \right)$$

where ties are broken arbitrary. Let $Q_Y^{\mathcal{C}}$ distribute uniformly over the codewords, i.e. $Q_Y^{\mathcal{C}}(y) = M^{-1}$ if $y \in \mathcal{C}$ and $Q_Y^{\mathcal{C}}(y) = 0$ otherwise. Then:

$$\tilde{W}_{Y|X=x} = W_{Y|X=x}^{M^{-1}}$$

Proof. Recall that

$$p_{c,x,y,U} = Q_Y^{\mathcal{C}} \left(d(x,Y) < d(x,y) \right) + U \cdot Q_Y^{\mathcal{C}} \left(d(x,Y) = d(x,y) \right)$$

 $\underbrace{ \text{Case I: } Q_Y^{\mathcal{C}}\left(d(x,Y) < d(x,y)\right) > 0 : \text{ In this case there exist } y' \neq y \text{ such that } d(x,y') < d(x,y). \text{ Thus: } \tilde{W}_{Y|X=x}(y|x) = 0 \\ \text{since } x \text{ cannot encode to } y (y' \text{ produce lower distortion}). \text{ From } Q_Y^{\mathcal{C}}\left(d(x,Y) < d(x,y)\right) \geq M^{-1} \text{ it follow that } p_{c,x,y,U} > M^{-1} \\ \text{with probability 1 and } W_{Y|X=x}^{M^{-1}}(y|x) = 0 \text{ as well.} \\ \underline{\text{Case II: } Q_Y^{\mathcal{C}}\left(d(x,Y) < d(x,y)\right) = 0 : \text{ Let:} }$

$$Q_Y^{\mathcal{C}}\left(d(x,Y) = d(x,y)\right) = k \cdot M^{-1}$$

where k is integer greater than 0. There are k different symbols $y_1 = y, \ldots, y_k \in \mathcal{Y}$ such that $d(x, y_i) = d(x, y)$ and $W_{Y|X=x}(\cdot|x)$ encode x to one of these symbols randomly. Thus:

$$\tilde{W}_{Y|X=x}(y_1|x) = \dots = \tilde{W}_{Y|X=x}(y_k|x) = k^{-1}$$

in particular, $\tilde{W}_{Y|X=x}(y|x) = k^{-1}$. Since:

$$p_{c,x,y,U} = U \cdot Q_Y^{\mathcal{C}} \left(d(x,Y) = d(x,y) \right) = U \cdot k \cdot M^{-1}$$

we have:

$$\mathbb{P}_{U}\left\{p_{c,x,y,U} \le M^{-1}\right\} = \mathbb{P}_{U}\left\{U \cdot k \cdot M^{-1} \le M^{-1}\right\} = k^{-1}$$

thus:

$$W_{Y|X=x}^{M^{-1}}(y) = M \cdot Q_Y^{\mathcal{C}}(y) \cdot \mathbb{P}_U \left\{ p_{c,x,y,U} \le M^{-1} \right\} = k^{-1}$$

as required.

Theorem 7. For any code C:

$$D(P_X, \mathcal{C}) = \tilde{D}(M^{-1}, Q_V^{\mathcal{C}})$$

in particular:

$$D(P_X, R) \ge \inf_{Q_Y \in \mathscr{P}(\mathcal{Y})} \tilde{D}(e^{-R}, Q_Y)$$

Proof.

$$D(P_X, \mathcal{C}) = \mathbb{E}_{P_X \times \tilde{W}_{Y|X}} (d(X, Y))$$
$$= \mathbb{E}_{P_X \times W_{Y|X}^{M^{-1}}} (d(X, Y))$$
$$= \tilde{D}(M^{-1}, Q_Y^{\mathcal{C}})$$

Let:

$$\hat{D}(z) \triangleq \inf_{Q_Y} \tilde{D}(e^{-z}, Q_Y)$$

Writing the achievable and converse results in terms of $\hat{D}(z)$, we have:

$$\hat{D}(R) \le D(P_X, R)$$

$$\le \inf_{\lambda \le R} \left\{ \hat{D}(R - \lambda) + d_{max} \cdot e^{-e^{\lambda}} (e^{\lambda} + 1) \right\}$$
(10)

This equation exemplify the tightness of our bounds.

V. VARIATIONAL FORMS OF $\tilde{D}(z, Q_Y)$

In this section we present two alternative presentations of $\tilde{D}(z, Q_Y)$. We first recall some definition. The definition of $\beta_{\alpha}(P,Q)$ which represent the optimal performance of a binary hypothesis testing between two σ -finite measures P and Q over a set W:

$$\beta_{\alpha}(P,Q) = \min_{\substack{P_{Z|W}:\\\sum_{w \in W} P(w)P_{Z|W}(1|w) \ge \alpha}} \sum_{w \in W} Q(w) P_{Z|W}(1|w), \tag{11}$$

 $P_{Z|W}: W \to \{0,1\}$ is any randomized test between P and Q. The minimum is guaranteed to be achieved by the Neyman–Pearson lemma. The functional $\beta_{\alpha}(P,Q)$ has been proved useful for converse results in channel coding and lossy compression, *e.g.* [3, Theorem 26], [4, Theorem 8].

The ∞ -order divergences between P and Q is:

$$D_{\infty}(P||Q) \triangleq \log \inf \left\{ \lambda : P(x) \le \lambda Q(x), \forall x \right\}$$
(12)

Theorem 8 (Variational forms of $\tilde{D}(z, Q_Y)$). Let $w = e^{-R}$. Then:

$$\tilde{D}(z, Q_Y) = \sup_{\substack{Q_X \\ W_Y|_X: D_{\infty}(P_X \times W_Y|_X ||P_X \times Q_Y) \le R}} \beta_w \left(Q_X \times Q_Y, P_X \times Q_Y \times d(X, Y)\right)$$
(13)
(14)

Corollary 9. The convexity of $\tilde{D}(z, Q_Y)$ with respect to Q_Y readily follows from (13). Since $\beta_w (Q_X \times Q_Y, P_X \times Q_Y \times d(X, Y))$ is convex with respect to Q_Y according to [1, Theorem 6] and supremum of convex function is convex as well.

Proof. To prove (13), let $\mathbb{1}_{\{p_{c,X,Y,U} \leq w\}}$ denote a (not necessarily optimal) test between $Q_X \times Q_Y$ and $P_X \times Q_Y \times d(X,Y)$. Since $\mathbb{P}_{Q_Y} \{p_{c,X,Y,U} \leq w\} = w$ for any x, it follow that for any Q_X we have: $\mathbb{P}_{Q_X \times Q_Y} \{p_{c,X,Y,U} \leq w\} = w$:

$$\beta_w \left(Q_X \times Q_Y, P_X \times Q_Y \times d(X, Y) \right) \\ \leq \mathbb{E}_{P_X \times Q_Y \times d(X, Y)} \left(\mathbb{1}_{\{p_{c, X, Y, U} \le w\}} \right) \\ = \mathbb{E}_{P_X \times Q_Y} \left(d(X, Y) \cdot \mathbb{1}_{\{p_{c, X, Y, U} \le w\}} \right)$$

Thus:

$$\sup_{Q_X} \beta_u \left(Q_X \times Q_Y, P_X \times Q_Y \times d(X, Y) \right) \\ \leq \mathbb{E}_{P_X \times Q_Y} \left(d(X, Y) \cdot \mathbb{1}_{\{p_{c,X,Y,U} \le u\}} \right)$$

To show the reverse inequality, we show that there exist Q_X^c such that:

$$\beta_w \left(Q_X^c \times Q_Y, P_X \times Q_Y \times d(X, Y) \right) \\= \mathbb{E}_{P_X \times Q_Y} \left(d(X, Y) \cdot \mathbb{1}_{\{p_{c,X,Y,U} \le w\}} \right)$$

which will complete the proof. We will construct such a Q_X^c and show that the optimal likelihood test between $Q_X^c \times Q_Y$ and $P_X \times Q_Y \times d(X,Y)$ matches the test $\mathbb{1}_{\{p_{c,X,Y,U} \leq w\}}$. The likelihood ratio is:

$$L(x,y) = \frac{P_X(x) \times Q_Y(y) \cdot d(x,y)}{Q_X^c(x) \times Q_Y(y)} = \frac{P_X(x)}{Q_X^c(x)} \times d(x,y)$$

and the likelihood ratio test is:

$$P(z|x,y) = \mathbb{1}_{\{L(x,y) < \lambda\}} + \sum_{x'} \tau_{x'} \mathbb{1}_{\{L(x,y) = \lambda, x = x'\}}$$

where λ and $\tau_{x'}$ are tuned so that:

$$\mathbb{P}_{Q_X^c \times Q_Y} \left\{ P(z|X,Y) = 1 \right\} = w$$

For any x there exist y_x and τ_x such that:

$$w = Q_Y(d(x, Y) < d(x, y_x)) + \tau_x \cdot Q_Y(d(x, Y) = d(x, y_x))$$

Define λ and Q_X^c :

$$\lambda = \sum_{x} P_X(x) \cdot d(x, y_x)$$
$$Q_X^c(x) = \lambda^{-1} \cdot P_X(x) \cdot d(x, y_x)$$

Note that Q_X^c is probability distribution and:

$$L(x, y_x) = \frac{P_X(x)}{Q_X^c(x)} \cdot d(x, y_x) = \lambda$$

To show that the tests matches we have to prove:

$$\mathbb{P}\left\{p_{c,x,y,U} \le w\right\} = \mathbb{P}\left\{P(z|x,y) = 1\right\}$$

for any x and y. There are 3 cases to consider:

1) $d(x,y) < d(x,y_x)$: In this case:

$$p_{c,x,y,U} \le Q_Y(d(x,Y) \le d(x,y))$$
$$\le Q_Y(d(x,Y) < d(x,y_x))$$
$$< w$$

thus: $\mathbb{P}\left\{p_{c,x,y,U} \leq w\right\} = 1$. On the other hand, since:

$$L(x,y) = \frac{P_X(x)}{Q_X^c(x)} \cdot d(x,y)$$
$$< \frac{P_X(x)}{Q_X^c(x)} \cdot d(x,y_x)$$
$$= \lambda$$

we also have: $\mathbb{P}\left\{P(z|x,y)=1\right\}=1$. 2) $d(x,y) > d(x,y_x)$: In this case:

$$p_{c,x,y,U} \ge Q_Y(d(x,Y) < d(x,y))$$
$$\ge Q_Y(d(x,Y) \le d(x,y_x))$$
$$\ge w$$

thus: $\mathbb{P} \{ p_{c,x,y,U} \leq w \} = 0$. On the other hand, since

$$L(x,y) = \frac{P_X(x)}{Q_X^c(x)} \cdot d(x,y)$$
$$> \frac{P_X(x)}{Q_X^c(x)} \cdot d(x,y_x)$$
$$= \lambda$$

we also have: $\mathbb{P} \{ P(z|x, y) = 1 \} = 0.$ 3) $d(x, y) = d(x, y_x)$: In this case

$$\mathbb{P}\left\{p_{c,x,y,U} \le w\right\} = \mathbb{P}\left\{p_{c,x,y_x,U} \le w\right\}$$
$$= \tau_x$$

On the other hand, since

$$L(x,y) = \frac{P_X(x)}{Q_X^c(x)} \cdot d(x,y)$$
$$= \frac{P_X(x)}{Q_X^c(x)} \cdot d(x,y_x)$$
$$= \lambda$$

we also have: $\mathbb{P}\left\{P(z|x,y)=1\right\}=\tau_x.$

Thus, the test $\mathbb{1}_{\{p_{c,x,y,U} \leq w\}}$ matches the likelihood ratio test and is, in fact, optimal. To prove (14), recall the channel:

$$W_{Y|X=x}^{e^{-R}} = e^{R} \cdot Q_{Y} \cdot \mathbb{1}_{\{p_{c,x,Y,U} \le e^{-R}\}}$$

For any y such that $Q_Y(y) > 0$:

$$\log \frac{e^R \cdot Q_Y(y) \cdot \mathbb{P}\left\{p_{c,x,Y,U} \le e^{-R}\right\}}{Q_Y(y)} \le R$$

Thus: $D_{\infty}(P_X \times W_{Y|X}^{e^{-R}} || P_X \times Q_Y) \leq R$ and:

 $\tilde{D}(e^{-R}, Q_Y) \geq \min_{\substack{W_{Y|X}: D_{\infty}(P_X \times W_{Y|X} || P_X \times Q_Y) \le R}} \mathbb{E}_{P_X \times W_{Y|X}} \left(d(X, Y) \right)$

On the other hand, if $W_{Y|X}$ satisfy:

 $D_{\infty}(P_X \times W_{Y|X} || P_X \times Q_Y) \le R$

we have $W_{Y|X}(y|x)P_X(x) \leq e^R \cdot Q_Y(y)P_X(x)$ for each x and y. Thus, to get the minimal $\mathbb{E}_{P_X \times W_{Y|X}}(d(X,Y))$ we will have to assign the maximal probability to the minimal distortion, *i.e.* $Q_Y(y) \cdot e^R$. This is exactly what the assignment $Q_Y \cdot \mathbb{1}_{\{p_{c,X,Y,U} \leq e^{-R}\}} \cdot e^R$ does which assign $Q_Y(y) \cdot e^R$ for the smallest distortion values until it exhaust the probability to one.

Let $W_{Y|X}$ denote any channel and let Q_Y denote the marginal distribution of $W_{Y|X} \times P_X$. Let $i(x; y) = \frac{W_{Y|X}(y|x)}{Q_Y(y)}$ denote the information density. Note that in theorem 8 we did not require that Q_Y is the marginal distribution. We might have relaxed the requirement $D_{\infty}(P_X \times W_{Y|X} || P_X \times Q_Y) \leq R$ which amounts to $i(x; y) \leq R$ for each x and y when Q_Y is the marginal distribution to the requirement that $\mathbb{P}_{P_X \times W_{Y|X}} \{i(X, Y) \leq R - \delta\}$ is close to 1.

Theorem 10. Let $W_{Y|X}$ such that:

$$\mathbb{P}_{P_X \times W_{Y|X}} \{ i(X, Y) \le R - \delta \} = e^{-\lambda}$$

Then:

$$\tilde{D}(e^{-(R-\delta)-\lambda}, Q_Y) \le \mathbb{E}_{P_X \times W_{Y|X}}(d(X, Y)) \cdot e^{\frac{1}{2}}$$

Proof.

$$\mathbb{E}_{P_X \times W_{Y|X}} \left(d(X,Y) \right) \ge \mathbb{E}_{P_X \times W_{Y|X}} \left(d(X,Y) \mathbb{1}_{\{i(x,y) \le R-\delta\}} \right)$$
$$\ge \mathbb{E}_{P_X \times W'_{Y|X}} \left(d(X,Y) \right) e^{-\lambda}$$
$$\stackrel{(a)}{\ge} \tilde{D}(e^{-(R-\delta)-\lambda}, Q_Y) \cdot e^{-\lambda}$$

where $W'_{Y|X} = e^{\lambda} \cdot W_{Y|X} \mathbb{1}_{\{i(x,y) \leq R-\delta\}}$ is a probability distribution. Note that:

$$\frac{W'_{Y|X}(y|x)}{Q_Y(y)} = \frac{e^{\lambda} \cdot W_{Y|X}(y|x) \mathbb{1}_{\{i(x,y) \le R-\delta\}}}{Q_Y(y)} \le e^{R-\delta} \cdot e^{\lambda}$$

Thus $D_{\infty}(P_X \times W'_{Y|X} || P_X \times Q_Y) \le R - \delta + \lambda$ and (a) follow from (14).

VI. EXCESS DISTORTION

The excess distortion is a spacial case of the average distortion that we have analyzed. Let d_{th} denote the target distortion level, replacing d(x, y) with $\mathbb{1}_{\{d(x,y)>d_{th}\}}$. Let:

$$\hat{D}(R, d_{th}, Q_Y) = \mathbb{E}_{P_X \times Q_Y} \left(\mathbb{1}_{\{d(X, Y) > d_{th}\}} \mathbb{1}_{\{p_{c, X, Y, U} < e^{-R}\}} \right)$$
(15)

Note that $p_{c,x,y,u}$ is also defined with respect to the "new" distortion:

$$p_{c,x,y,u} = Q_Y \left(\mathbb{1}_{\{d(x,Y) > d_{th}\}} < \mathbb{1}_{\{d(x,y) > d_{th}\}} \right) + u \cdot Q_Y \left(\mathbb{1}_{\{d(x,Y) > d_{th}\}} < \mathbb{1}_{\{d(x,y) > d_{th}\}} \right)$$

equation (14) translates in this case to:

$$\dot{D}(R, d_{th}, Q_Y) = \inf_{\substack{W_Y|X:\\D_{\infty}(P_X \times W_Y|X} ||P_X \times Q_Y) \le R}} \mathbb{P}_{P_X \times W_Y|X} \left\{ d(X, Y) > d_{th} \right\}$$

and (8) is:

$$\tilde{R}(\delta, z, Q_Y) = \inf_{\substack{W_Y|_X:\\ \mathbb{P}_{P_X \times W_Y|_X} \{d(X, Y) > z\} \le \delta}} D_{\infty}(P_X \times W_Y|_X ||P_X \times Q_Y)$$

VII. RELATION TO KNOWN BOUNDS

The information spectrum approach [5, Theorem 5.5.1] provides a general formula for the rate distortion function. The general formula takes any channel $W_{Y|X}$ and with slight abuse of notation, say that the distortion $\mathbb{E}_{P_X \times W_{Y|X}}(d(X,Y))$ is achievable for a code with rate R such that $\mathbb{P}_{P_X \times W_{Y|X}} \{i(x;y) \leq R\}$ approach 1 where i(x;y) is the information density. The achievability is proved by the random coding argument where the distribution used to draw the codewords is the marginal distribution of $P_X \times W_{Y|X}$ on \mathcal{Y} . In this paper (see also [6]) we started with any distribution Q_Y and analyzed the random code performance with no channel in mind. For the achievable part, theorem 10 provides the link between the functional $\tilde{D}(z, Q_Y)$ and the elements in the information spectrum formula. The converse follow since a code \mathcal{C} with rate R satisfies

$$D_{\infty}\left(P_X \times \tilde{W}_{Y|X} || P_X \times Q_Y^{\mathcal{C}}\right) = R$$

where $\tilde{W}_{Y|X}$ and $Q_Y^{\mathcal{C}}$ were defined in the text.

Much attention has been given to the problem of lossy compression with the excess distortion constraint. The tightest results (to the best of our knowledge) appeared in [6]. The converse bound [6, Theorem 2] was shown to be tighter than [7, Theorem 8]. In [8] the author demonstrated how the bound [7, Theorem 8] can be relaxed to all other bounds presented there.

The approach in [6] for the achievability was to analyze the random coding for a given prior distribution Q_Y and bound the performance from above using a change of measure from $P_X \times Q_Y$ to $P_X \times W_{Y|X}$ where the marginal distribution of $P_X \times W_{Y|X}$ with respect to \mathcal{Y} does not necessarily matches Q_Y . Later they optimize for a given channel $W_{Y|X}$ the best prior distribution. Thus, their bounds are given as an optimization over the a set of channels. In [6, Theorem 2] the author defined:

$$M(P_{XY}) \triangleq \sum_{y \in \mathcal{Y}} \sup_{x \in \mathcal{X}: P_{XY}(x,y) > 0} P_{Y|X}(y|x)$$

and [6, Lemma 4] reads:

$$\inf_{Q_Y \in \mathscr{P}(\mathcal{Y})} D_{\infty}(P_{XY} || P_X \times Q_Y) = \log M(P_{XY})$$

Hence:

$$\inf_{\substack{Q_Y \in \mathscr{P}(\mathcal{Y}) \\ P_{Y \times W_{Y|X}}: \\ \mathbb{P}_{P_X \times W_{Y|X}} \{d(X,Y) > d_{th}\} \le \delta}} \log M(P_X \times W_{Y|X})$$

and our converse bound matches theirs. Their achievability bound (Theorem 2) is slightly tighter than the bound in corollary 4. While their bound reads:

$$R \le \min_{\delta < \epsilon} \tilde{R}(\delta, d_{th}) + \log \log \left(\frac{1 - \delta}{\epsilon - \delta}\right)$$

our bound is:

$$R \le \min_{\delta < \epsilon} \tilde{R}(\delta, d_{th}) + g\left(\frac{1-\delta}{\epsilon - \delta}\right)$$

where:

$$g(x) = \log \log(x) + \log\left(\frac{2}{3}\right) + \log\left(1 + \sqrt{1 + 9\left(2\log(x)\right)^{-1}}\right)$$

The difference smaller than 1 nat for all practical cases. Note that, our bound follow from bounding

$$\tilde{D}(z,Q_Y) \le \tilde{D}(u,Q_Y) \mathbb{1}_{\{z \le u\}} + \tilde{D}(1,Q_Y) \mathbb{1}_{\{z > u\}}$$

and evaluating the exact formula with this bound. Using the exact formula directly or bound $\tilde{D}(z, Q_Y)$ using more points would gives tighter bound. Since their achievability also follow from the random coding, obviously the exact formula is tighter.

VIII. PRIOR OPTIMIZATION

The prior optimization problem is the main drawback of our approach. While the optimization problem is convex, we do not have a close form solution to the important case of memoryless source (and channel in the channel coding case) and this should be further investigated. In the channel coding case, the prior optimization problem for memoryless channel is "solved" by [9, Theorem 10] which shows that for memoryless channel, we can resort to memoryless priors. For our functional $\tilde{D}(z, Q_Y)$ (and the one used in the channel coding problem) we do not have such a theorem. Moreover, simulation results show that this is not even true, and sometimes, prior with memory are better than the memoryless priors.

IX. CONCLUDING REMARKS

In this paper we presented a novel analysis for the lossy compression problem. We have analyzed the general case of average distortion constraints. We presented tight achievable and converse bounds. Both bounds are given in terms of the functional $\tilde{D}(z, Q_Y)$ which has resemblance to the meta-converse in channel coding. The problem of finding distribution minimizing $\tilde{D}(z, Q_Y)$ is still open in general although the $\tilde{D}(z, Q_Y)$ is convex with respect to Q_Y .

APPENDIX

Lemma 1. For any $x \in (0,1)$ let λ be the solution to $x = e^{-e^{\lambda}}(e^{\lambda} + 1)$, then:

$$\lambda - \log(-\log(x)) \le \log(\frac{2}{3}) + \log\left(1 + \sqrt{1 - \frac{9}{2\log(x)}}\right)$$
 (16)

$$\lambda - \log(-\log(x)) \ge -\log(2) + \log\left(1 + \sqrt{1 - \frac{8}{\log(x)}}\right) \tag{17}$$

Proof. Let $t = e^{\lambda}$ and z = -log(x). Then:

$$z = t - \log(1+t)$$

Using: $\log(1+t) \le \frac{t(6+t)}{2(3+2t)}$ ([10, Eq. 22]) we have:

$$z \ge t - \frac{t(6+t)}{2(3+2t)} = \frac{t^2}{2 + \frac{4}{3}t}$$

For $t, z \ge 0$ the only solution to: $z = \frac{t^2}{2+\frac{3}{3}t}$ is: $t = \frac{2}{3}z\left(1 + \sqrt{1+\frac{9}{2z}}\right)$. Thus: $z \ge \frac{3t^2}{2(3+2t)}; t, z \ge 0 \iff$ $t \le \frac{2}{3}z\left(1 + \sqrt{1+\frac{9}{2z}}\right)$

Hence:

$$\lambda = \log(t)$$

$$\leq \log\left(\frac{2}{3}z\left(1 + \sqrt{1 + \frac{9}{2z}}\right)\right)$$

$$= \log(-\log(x)) + \log(\frac{2}{3}) + \log\left(1 + \sqrt{1 - \frac{9}{2\log(x)}}\right)$$

The lower bound follow along the same line using $\log(1+t) \ge \frac{2t}{2+t}$. **Proposition 11.** Let U_i denote M independent uniform random variables and let $W_M = \min \{U_i\}$. The c.d.f of W_M is:

$$f_M(w) = M \cdot (1-u)^{M-1}$$
(18)

Proof.

$$\mathbb{P}\left\{W_M \le w\right\} = 1 - \mathbb{P}\left\{\min_i \left\{U_i\right\} > w\right\}$$
$$= 1 - \prod_{i=1}^M \mathbb{P}\left\{U_i > w\right\}$$
$$= 1 - (1 - w)^M$$

Thus, the p.d.f. of W_M is $f_M(w) = M \cdot (1-u)^{M-1}$.

REFERENCES

- [1] Y. Polyanskiy, "Saddle point in the minimax converse for channel coding," *Information Theory, IEEE Transactions on*, vol. 59, no. 5, pp. 2576–2595, 2013.
- [2] N. Elkayam and M. Feder, "Achievable and converse bounds over a general channel and general decoding metric," arXiv preprint arXiv:1411.0319, 2014. [Online]. Available: http://www.eng.tau.ac.il/~elkayam/FiniteBlockLen.pdf
- [3] Y. Polyanskiy, H. V. Poor, and S. Verdú, "Channel coding rate in the finite blocklength regime," *Information Theory, IEEE Transactions on*, vol. 56, no. 5, pp. 2307–2359, 2010.
- [4] V. Kostina and S. Verdú, "Lossy joint source-channel coding in the finite blocklength regime," *IEEE Transactions on Information Theory*, vol. 59, no. 5, pp. 2545–2575, 2013.
- [5] H. Koga et al., Information-spectrum methods in information theory. Springer, 2002, vol. 50.
- [6] T. Matsuta and T. Uyematsu, "Non-asymptotic bounds for fixed-length lossy compression," in 2015 IEEE International Symposium on Information Theory (ISIT). IEEE, 2015, pp. 1811–1815.
- [7] V. Kostina and S. Verdú, "Fixed-length lossy compression in the finite blocklength regime," *Information Theory, IEEE Transactions on*, vol. 58, no. 6, pp. 3309–3338, 2012.
- [8] L. Palzer and R. Timo, "A converse for lossy source coding in the finite blocklength regime," in 24th International Zurich Seminar on Communications (IZS). ETH-Zürich, 2016.
- [9] S. Verdú and T. S. Han, "A general formula for channel capacity," IEEE Transactions on Information Theory, vol. 40, no. 4, pp. 1147–1157, 1994.
- [10] F. Topsok, "Some bounds for the logarithmic function," Inequality theory and applications, vol. 4, p. 137.