

Error exponents for source coding under logarithmic loss

Citation for published version (APA):

Joudeh, H., & Wu, H. (2024). Error exponents for source coding under logarithmic loss. In A. Lapidoth, & S. M. Moser (Eds.), *International Zurich Seminar on Information and Communication (IZS 2024): Proceedings* (pp. 99-103). ETH Zürich. <https://doi.org/10.3929/ethz-b-000664598>

DOI:

[10.3929/ethz-b-000664598](https://doi.org/10.3929/ethz-b-000664598)

Document status and date:

Published: 06/03/2024

Document Version:

Accepted manuscript including changes made at the peer-review stage

Please check the document version of this publication:

- A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.
- The final author version and the galley proof are versions of the publication after peer review.
- The final published version features the final layout of the paper including the volume, issue and page numbers.

[Link to publication](#)

General rights

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

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

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

Error exponents for source coding under logarithmic loss

Hamdi Joudeh and Han Wu

ICT Lab, Eindhoven University of Technology, The Netherlands

Abstract—In source coding under the logarithmic loss distortion measure, a source is compressed into a message, which is then decompressed into a soft reconstruction (i.e. probability distribution). The distortion is measured by the remaining uncertainty about the source given the message. Shkel and Verdú showed that this lossy source coding setting is intimately related to almost lossless source coding with list decoding, and used this insight to characterize the single-shot excess distortion error probability. In this work, we build upon this connection to list decoding and derive error exponents for source coding under logarithmic loss, without and with side information. The error exponents are closely related to their almost lossless counterparts.

I. INTRODUCTION

In this paper we study the problem of fixed-length lossy source coding of a discrete memoryless source (DMS) under the logarithmic loss (log-loss) distortion measure. While the log-loss is most commonly used in prediction and learning theory, its adoption as a distortion measure in lossy source coding is also natural, specifically in settings where the decoder produces a soft reconstruction (i.e. probability distribution) of the source instead of a point estimate [1]–[3].

The log-loss distortion measure enjoys some mathematical properties that enable elegant characterizations in a number of settings. For instance, under an average distortion criterion, the rate-distortion function is given by [1, Example 2]

$$R(\Delta) = H(X) - \Delta \quad (1)$$

where $H(X)$ is the source entropy and Δ is the average log-loss distortion (assume $0 \leq \Delta \leq H(X)$). The converse for the corresponding coding theorem is obtained by bounding the average log-loss distortion using the conditional entropy of the source given its reconstruction. By building upon this property, Courtade and Weissman [2] derived tight outer bounds in various multi-terminal source coding settings under average log-loss distortion (see also Courtade and Wesel [1]).

More recently, Shkel and Verdú [3] derived single-shot bounds under both excess and average log-loss distortion criteria, without and with decoder side information (see [4], [5] for universal extensions). Key to their approach is a close connection between the log-loss setting and the almost lossless setting with list decoding. As we shall see, this connection to list decoding also plays a central role in our current work.

In this paper, instead of single-shot bounds, we are interested in error exponents under an excess log-loss distortion criterion. We derive error exponents without and with side information,

while mainly focusing on universal schemes. We also demonstrate close connections to results in almost lossless settings.

Notation: We use standard notation, briefly explained here. $\mathcal{P}(\mathcal{X})$ denotes the probability simplex on a finite alphabet \mathcal{X} . For a probability mass function (pmf) $P_X \in \mathcal{P}(\mathcal{X})$, we denote its support by $\mathcal{S}(P_X)$ and its entropy by $H(P_X)$. The relative entropy between two pmfs Q_X and P_X is denoted by $D(Q_X \| P_X)$. For (X, Y) with joint pmf $P_{XY} = P_{X|Y}P_Y$, the conditional entropy of X given Y is denoted by $H(P_{X|Y}|P_Y)$. The type of a sequence $\mathbf{x} \in \mathcal{X}^n$ is denoted by $P_{\mathbf{x}}$, and $\mathcal{P}_n(\mathcal{X})$ is the set of all types of sequences in \mathcal{X}^n . For $Q \in \mathcal{P}_n(\mathcal{X})$, the corresponding type class is denoted by $\mathcal{T}_n(Q)$. Given a second sequence $\mathbf{y} \in \mathcal{Y}^n$, $P_{\mathbf{x}\mathbf{y}}$ and $P_{\mathbf{x}|\mathbf{y}}$ are the joint and conditional types. $\mathcal{P}_n(\mathcal{X} \times \mathcal{Y})$ and $\mathcal{P}_n(\mathcal{X}|\mathcal{Y})$ are the sets of all joint and conditional types. $\mathcal{T}_n(Q_{X|Y}|\mathbf{y})$ is the conditional type class of $Q_{X|Y} \in \mathcal{P}_n(\mathcal{X}|\mathcal{Y})$ given \mathbf{y} . We will make use of types and type classes and their basic properties, such as cardinality and probability bounds (see, e.g., [6, Ch.2]).

II. SOURCE CODING UNDER LOG-LOSS

Consider a DMS with finite alphabet \mathcal{X} that randomly generates i.i.d. source sequences $\mathbf{X} \triangleq (X_1, X_2, \dots, X_n)$ according to a pmf $P_X \in \mathcal{P}(\mathcal{X})$. We use $\mathbf{x} \triangleq (x_1, x_2, \dots, x_n)$ to denote a realization of \mathbf{X} . A *soft reconstruction* of \mathbf{x} is a member of $\mathcal{P}(\mathcal{X}^n)$, i.e. a distribution on \mathcal{X}^n , denoted by \hat{P}_n . The log-loss distortion between \mathbf{x} and \hat{P}_n is defined as

$$d(\mathbf{x}, \hat{P}_n) \triangleq \log \frac{1}{\hat{P}_n(\mathbf{x})}. \quad (2)$$

The log-loss, also referred to as the self-information loss, can be understood as the remaining uncertainty about \mathbf{x} given its reconstruction \hat{P}_n [1]–[3]. For instance, $d(\mathbf{x}, \hat{P}_n)$ is zero if and only if \hat{P}_n has a single mass point at \mathbf{x} , i.e. an exact *hard reconstruction*; and infinite whenever \mathbf{x} has zero probability under \hat{P}_n . For convenience, we work with the normalized (per-symbol) log-loss defined as $d_n(\mathbf{x}, \hat{P}_n) \triangleq \frac{1}{n}d(\mathbf{x}, \hat{P}_n)$.

In the lossy source coding setting considered in this work, the sequence \mathbf{X} is encoded into a message index from the finite set \mathcal{M}_n , which is then decoded into a soft reconstruction from $\mathcal{P}(\mathcal{X}^n)$. A lossy source code of block-length n is thus a pair of mappings $\phi_n : \mathcal{X}^n \rightarrow \mathcal{M}_n$ and $\varphi_n : \mathcal{M}_n \rightarrow \mathcal{P}(\mathcal{X}^n)$, referred to as the encoder and decoder respectively.

For a lossy source code (ϕ_n, φ_n) , the code *rate* is given by $\frac{1}{n} \log |\mathcal{M}_n|$, while $\mathbb{P} [d_n(\mathbf{X}, \varphi_n(\phi_n(\mathbf{X}))) > \Delta]$ is the excess distortion *error probability* for some distortion level $\Delta \geq 0$. We say that (ϕ_n, φ_n) is an (n, R, Δ, ϵ) -code if

$$\frac{1}{n} \log |\mathcal{M}_n| \leq R \quad \text{and} \quad \mathbb{P} [d_n(\mathbf{X}, \varphi_n(\phi_n(\mathbf{X}))) > \Delta] \leq \epsilon.$$

This work was partially supported by the European Research Council (ERC) under the ERC Starting Grant N. 101116550 (IT-JCAS).

The minimal error probability for fixed (n, R, Δ) is defined as

$$\varepsilon(n, R, \Delta) \triangleq \inf \{ \epsilon : \text{there exists an } (n, R, \Delta, \epsilon)\text{-code} \}.$$

We are interested in characterizing the asymptotic behaviour of $\varepsilon(n, R, \Delta)$, captured through the error exponent defined as

$$E(R, \Delta) \triangleq \lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{\varepsilon(n, R, \Delta)}. \quad (3)$$

Remark 1. The above problem does not fall under the umbrella of standard lossy source coding in discrete memoryless settings [6], [7]. In the standard paradigm, the reconstruction is a sequence drawn from a discrete product alphabet; and the distortion is additive, i.e. a normalized sum of single-letter distortions. In the setting considered here, the reconstruction alphabet $\mathcal{P}(\mathcal{X}^n)$ is not a product alphabet and is not discrete; and the distortion measure is not additive. For $d_n(\mathbf{x}, \hat{P}_n)$ to be additive, the soft reconstruction \hat{P}_n must be a product distribution, as in earlier works on log-loss source coding [1], [2]. This need not be the case in general, and as we shall see, the codes we propose employ non-product soft reconstructions.

A. Connection to list decoding

In [3], Shkel and Verdú established a fundamental connection between lossy source coding under log-loss and almost lossless source coding with list decoding, leading to an exact characterization of $\varepsilon(n, R, \Delta)$. This connection is central to the approach we take here, therefore, we review it in some detail. We start with a key lemma linking the log-loss of a soft reconstruction to the list size in list decoding.

To this end, fix a soft reconstruction $\hat{P}_n \in \mathcal{P}(\mathcal{X}^n)$ and a distortion level $\Delta \geq 0$. We say that a sequence $\mathbf{x} \in \mathcal{X}^n$ is Δ -covered by \hat{P}_n if $d_n(\mathbf{x}, \hat{P}_n) \leq \Delta$. If \hat{P}_n Δ -covers every element of a set (or list) $\mathcal{L}_n \subseteq \mathcal{X}^n$, then the set \mathcal{L}_n is also said to be Δ -covered by the soft reconstruction \hat{P}_n .

Lemma 1. *Let $\mathcal{L}_n \subseteq \mathcal{X}^n$. There exists a soft-reconstruction $\hat{P}_n \in \mathcal{P}(\mathcal{X}^n)$ that Δ -covers \mathcal{L}_n if and only if*

$$|\mathcal{L}_n| \leq \lfloor \exp(n\Delta) \rfloor. \quad (4)$$

Proof. The direct part holds by taking \hat{P}_n to be uniform on \mathcal{L}_n and zero elsewhere. The converse part follows from [3, Lemma 1], reproduced here for completeness. Let $\mathcal{B}_n(\Delta, \hat{P}_n)$ be the set of all source sequences Δ -covered by \hat{P}_n , i.e.

$$\mathcal{B}_n(\Delta, \hat{P}_n) \triangleq \left\{ \mathbf{x} \in \mathcal{X}^n : d_n(\mathbf{x}, \hat{P}_n) \leq \Delta \right\}. \quad (5)$$

It is sufficient to show $|\mathcal{B}_n(\Delta, \hat{P}_n)| \leq \lfloor \exp(n\Delta) \rfloor$. Note that $\mathbf{x} \in \mathcal{B}_n(\Delta, \hat{P}_n)$ implies $\hat{P}_n(\mathbf{x}) \geq \exp(-n\Delta)$, and therefore

$$1 = \sum_{\mathbf{x} \in \mathcal{X}^n} \hat{P}_n(\mathbf{x}) \geq \sum_{\mathbf{x} \in \mathcal{B}_n(\Delta, \hat{P}_n)} \hat{P}_n(\mathbf{x}) \geq |\mathcal{B}_n(\Delta, \hat{P}_n)| \exp(-n\Delta).$$

The bound is tightened by including the floor function. \square

Now consider an almost lossless list source code: here a source sequence is encoded into one of $\lfloor \exp(nR) \rfloor$ message indices, and a message index is decoded into a list of $\lfloor \exp(n\Delta) \rfloor$ source sequences. A decoding error occurs if the generated

source sequence is not in the decoded list. From Lemma 1, we see that such a code with error probability ϵ can be converted into a (n, R, Δ, ϵ) log-loss source code. Conversely, a log-loss source code can be converted into an almost lossless list source code. This connection leads to the following result.

Theorem. (Shkel-Verdú [3, Theorem 5-6]). *Let $G : \mathcal{X}^n \rightarrow \{1, 2, \dots, |\mathcal{X}^n|\}$ be a probability rank function that ranks source sequences in decreasing order of their probability. Then*

$$\varepsilon(n, R, \Delta) = \mathbb{P}[G(\mathbf{X}) > \lfloor \exp(nR) \rfloor \lfloor \exp(n\Delta) \rfloor]. \quad (6)$$

III. ERROR EXPONENT

We now present the first result of this paper. To this end, we first define the function $F(R)$ on $0 \leq R < \log |\mathcal{S}(P_X)|$ as

$$F(R) \triangleq \min_{Q_X : H(Q_X) \geq R} D(Q_X \| P_X). \quad (7)$$

Note that $F(R) = E(R, 0)$, which is the error exponent in the almost lossless case [6], [8], [9]. $F(R)$ is continuous, convex and increasing on its domain, with $F(R) = 0$ on $0 \leq R \leq H(P_X)$. The expression in (7) is known as a *primal form*. $F(R)$ admits an equivalent *dual form* given as [6, Prob.2.15]

$$F(R) = \sup_{\rho \geq 0} \rho \left(R - H_{\frac{1}{1+\rho}}(X) \right) \quad (8)$$

where $H_{\frac{1}{1+\rho}}(X)$ is the Rényi entropy of order $1/(1+\rho)$.

Recall that the log-loss rate-distortion function is given by $R(\Delta) = H(P_X) - \Delta$, and to ensure that $\varepsilon(n, R, \Delta)$ goes to zero, we must have $R > H(P_X) - \Delta$. On the other hand, for rates satisfying $R \geq \log |\mathcal{S}(P_X)| - \Delta$, the whole support of P_X^n can be covered by lists of size $e^{n\Delta}$. Therefore, the relevant range is $H(P_X) < R + \Delta < \log |\mathcal{S}(P_X)|$.

Theorem 1. *Let (R, Δ) be a rate-distortion pair such that $H(P_X) < R + \Delta < \log |\mathcal{S}(P_X)|$. Then*

$$E(R, \Delta) = F(R + \Delta). \quad (9)$$

We observe that $E(R, \Delta) = E(R + \Delta, 0)$, i.e. the log-loss error exponent as a function of R is merely a translation of the almost lossless error exponent by Δ bits to the left. This can be understood in light of the optimal source code that achieves (6) as follows. For a log-loss source code of rate R and distortion Δ , an excess distortion error occurs when the source generates a sequence with probability rank greater than $\lfloor e^{nR} \rfloor \lfloor e^{n\Delta} \rfloor$, that is the number of sequences covered by $\lfloor e^{nR} \rfloor$ lists of size $\lfloor e^{n\Delta} \rfloor$ each. On the other hand, an almost lossless source code of rate $R + \Delta$ makes an error when the source generates a sequence with probability rank greater than $\lfloor e^{n(R+\Delta)} \rfloor$. Asymptotically, the two error events have almost the same probability, yielding in the same error exponent.

The above argument is sufficient for proving Theorem 1, yet it employs a code that depends on the source pmf (or probability rank function). Further on we present an alternative proof using the method of types, extending the Longo-Sgarro approach [9] (see also Csiszár-Körner [6]) from the almost lossless case to the log-loss case. As is often the case with types-based proofs, a universal scheme emerges.

It is worthwhile mentioning that the log-loss error exponent expression in (9) can be obtained from Marton's error exponent [7] by replacing the general rate-distortion function with its log-loss counterpart. While this is perhaps expected, the result in Theorem 1 does not follow directly from Marton's proof, at least not without modification, as the log-loss setting considered here is not a special case of the classical rate-distortion setting (see Remark 1). We shall see next that Theorem 1 is proved directly using the connection to list decoding.

A. Proof of Theorem 1

1) *Achievability*: Fix $R > 0$ and $n \in \mathbb{N}$, and let $J_n = |\mathcal{P}_n(\mathcal{X})|$ and $M_n = \lfloor (1+n)^{-|\mathcal{X}|} e^{nR} \rfloor$. For every type $Q \in \mathcal{P}_n(\mathcal{X})$, partition $\mathcal{T}_n(Q)$ into M_n lists all roughly of the same size. $\mathcal{L}_n(\mathbf{x})$ denotes the list containing \mathbf{x} . By construction,

$$|\mathcal{L}_n(\mathbf{x})| \leq \left\lceil \frac{|\mathcal{T}_n(P_{\mathbf{x}})|}{M_n} \right\rceil \quad (10)$$

for every $\mathbf{x} \in \mathcal{X}^n$. An emitted source sequence \mathbf{x} is encoded into $\phi_n(\mathbf{x}) = (t_n(\mathbf{x}), l_n(\mathbf{x}))$, where $t_n(\mathbf{x}) \in \{1, 2, \dots, J_n\}$ is the type index and $l_n(\mathbf{x}) \in \{1, 2, \dots, M_n\}$ is the list index. Upon receiving $\phi_n(\mathbf{x})$, the decoder reproduces the list $\mathcal{L}_n(\mathbf{x})$ containing \mathbf{x} . The corresponding soft reconstruction $\hat{\varphi}_n(\phi_n(\mathbf{x})) = \hat{P}_n(\cdot | \phi_n(\mathbf{x}))$ is set as

$$\hat{P}_n(\hat{\mathbf{x}} | \phi_n(\mathbf{x})) = \begin{cases} \frac{1}{|\mathcal{L}_n(\mathbf{x})|}, & \hat{\mathbf{x}} \in \mathcal{L}_n(\mathbf{x}) \\ 0, & \text{otherwise.} \end{cases}$$

The rate of this code satisfies $\frac{1}{n} \log(J_n M_n) \leq R$. For any sequence $\mathbf{x} \in \mathcal{X}^n$, the log-loss incurred by the corresponding reconstruction $\varphi_n(\phi_n(\mathbf{x}))$ is bounded above as

$$\begin{aligned} d_n(\mathbf{x}, \varphi_n(\phi_n(\mathbf{x}))) &= \frac{1}{n} \log |\mathcal{L}_n(\mathbf{x})| \\ &\leq \frac{1}{n} \log \left[\frac{e^{nH(P_{\mathbf{x}})}}{\lfloor e^{nR - |\mathcal{X}| \log(1+n)} \rfloor} \right] \quad (11) \\ &\leq H(P_{\mathbf{x}}) - R + \delta_n \quad (12) \end{aligned}$$

where $\delta_n \geq 0$ goes to zero as n grows large.

We now analyze the error probability. To this end, fix $\Delta \geq 0$ such that $H(P_X) < R + \Delta < \log |\mathcal{S}(P_X)|$. We can see from (12) that all source sequences in the set \mathcal{B}_n , defined as

$$\mathcal{B}_n = \bigcup_{Q \in \mathcal{P}_n(\mathcal{X}): H(Q) \leq R + \Delta - \delta_n} \mathcal{T}_n(Q), \quad (13)$$

are reconstructed with a log-loss not exceeding Δ , and thus the excess distortion error event is included in $\mathcal{X}^n \setminus \mathcal{B}_n$. The error probability under source pmf P_X is bounded above as

$$\begin{aligned} P_X^n(\mathcal{X}^n \setminus \mathcal{B}_n) &= \sum_{Q \in \mathcal{P}_n(\mathcal{X}): H(Q) > R + \Delta - \delta_n} P_X^n(\mathcal{T}_n(Q)) \\ &\leq \sum_{Q \in \mathcal{P}_n(\mathcal{X}): H(Q) > R + \Delta - \delta_n} e^{-nD(Q \| P_X)} \quad (14) \end{aligned}$$

$$\leq (n+1)^{|\mathcal{X}|} \max_{Q \in \mathcal{P}_n(\mathcal{X}): H(Q) > R + \Delta - \delta_n} e^{-nD(Q \| P_X)} \quad (15)$$

$$\leq (n+1)^{|\mathcal{X}|} e^{-nF(R + \Delta - \delta_n)}. \quad (16)$$

The above steps are standard and use the properties of types and type classes. Achievability follows from (16) and the continuity of $F(R')$ on $H(P_X) < R' < \log |\mathcal{S}(P_X)|$.

Remark 2. As mentioned earlier, the above source code is universal and does not depend on P_X . The code is also universal with respect to the distortion level Δ , and only depends on the rate R (and block-length n). For fixed R , the same sequence of codes achieves a positive exponent for every P_X and Δ satisfying $R < H(P_X) - \Delta$. The key to the universality with respect to Δ is the *variable list partitioning* of type classes, where list sizes depend on the type and code rate but not on the distortion level (see (10)). The same partitioning is used by Bunte and Lapidoth in [10] in the context of strictly lossless list source coding (also known as task encoding), where the focus is on analyzing list size moments.

Remark 3. In the achievability proof of Marton's error exponent [7], a key ingredient is the type covering lemma which states that for any type $Q \in \mathcal{P}_n(\mathcal{X})$ with a rate-distortion function satisfying $R(Q, \Delta) \leq R - \delta$, the corresponding type class $\mathcal{T}_n(Q)$ can be Δ -covered by e^{nR} reconstruction sequences. The type covering lemma is often proved using random selection (i.e. random coding). In the log-loss setting considered here, type covering is accomplished through simple partitioning, and the proof does not rely on random coding.

2) *Converse*: Fix a pair (R, Δ) and a source pmf P_X such that $H(P_X) < R + \Delta < \log |\mathcal{S}(P_X)|$. For every block-length n , let (ϕ_n^*, φ_n^*) be an optimal code achieving the minimal error probability $\varepsilon(n, R, \Delta)$. Moreover, define the set $\mathcal{B}_n^* \triangleq \{\mathbf{x} \in \mathcal{X}^n : d_n(\mathbf{x}, \varphi_n^*(\phi_n^*(\mathbf{x}))) \leq \Delta\}$. An error occurs whenever the source produces a sequence in $\mathcal{X}^n \setminus \mathcal{B}_n^*$.

Let $M(\Delta, \mathcal{B}_n^*)$ be the minimum number of soft reconstructions required to Δ -cover \mathcal{B}_n^* . From Lemma 1, we know that any soft reconstruction can Δ -cover at most $\lfloor e^{n\Delta} \rfloor$ source sequence. Therefore, we must have

$$M(\Delta, \mathcal{B}_n^*) \geq \left\lceil \frac{|\mathcal{B}_n^*|}{\lfloor e^{n\Delta} \rfloor} \right\rceil. \quad (17)$$

It immediately follows that the rate R of (ϕ_n^*, φ_n^*) must satisfy

$$e^{nR} \geq M(\Delta, \mathcal{B}_n^*) \geq |\mathcal{B}_n^*| e^{-n\Delta}. \quad (18)$$

Now let $\delta_n = \frac{|\mathcal{X}|}{n} \log(1+n) + \frac{1}{n} \log 2$ and let $Q \in \mathcal{P}_n(\mathcal{X})$ be a type such that $H(Q) \geq R + \Delta + \delta_n$. The cardinality of the corresponding type class $\mathcal{T}_n(Q)$ is bounded below as

$$|\mathcal{T}_n(Q)| \geq (1+n)^{-|\mathcal{X}|} e^{nH(Q)} \geq 2e^{n(R+\Delta)} \geq 2|\mathcal{B}_n^*|$$

from which we conclude that at least half of the sequences in $\mathcal{T}_n(Q)$ are not contained in \mathcal{B}_n^* . Therefore

$$\begin{aligned} \varepsilon(n, R, \Delta) &= P_X^n(\mathcal{X}^n \setminus \mathcal{B}_n^*) \geq \frac{1}{2} P_X^n(\mathcal{T}_n(Q)) \\ &\geq \frac{1}{2} (n+1)^{-|\mathcal{X}|} e^{-nD(Q \| P_X)} = e^{-nD(Q \| P_X) - n\delta_n} \quad (19) \end{aligned}$$

obtained from the standard type class probability lower bound. This holds for all types satisfying $H(Q) \geq R + \Delta + \delta_n$, hence

$$-\frac{1}{n} \log \varepsilon(n, R, \Delta) \leq \delta_n + \min_{Q \in \mathcal{P}_n(\mathcal{X}): H(Q) \geq R + \Delta + \delta_n} D(Q \| P_X)$$

$$= \delta_n + F_n(R + \Delta + \delta_n) \quad (20)$$

where $F_n(R')$ is defined as $F(R')$ in (7) except that the minimization is over types in $\mathcal{P}_n(\mathcal{X})$ instead of all pmfs in $\mathcal{P}(\mathcal{X})$. By definition, we know that $F(R') \leq F_n(R')$. In addition, it can be shown that $F_n(R') \leq F(R') + \delta'_n$ for some $\delta'_n > 0$ that goes to zero as n grows large.¹ By combining this with (20) and taking the limit, the converse result follows.

IV. SIDE INFORMATION

In this section we consider settings with side information. Here we have a pair of DMSs with finite alphabets \mathcal{X} and \mathcal{Y} . The sources randomly generate an i.i.d. sequence of pairs $(\mathbf{X}, \mathbf{Y}) \triangleq ((X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n))$ according to a joint pmf $P_{XY} \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})$. The goal remains to compress the sequence \mathbf{X} and then decompress it into a soft reconstruction in $\mathcal{P}(\mathcal{X}^n)$, but now \mathbf{Y} is available either at both the encoder and decoder sides, or at the decoder side only.

A. Encoder-decoder side information

For the case where the side information sequence is available at both the encoder and decoder sides, a lossy source code of block-length n is given by the pair of mappings $\phi_n : \mathcal{X}^n \times \mathcal{Y}^n \rightarrow \mathcal{M}_n$ and $\varphi_n : \mathcal{M}_n \times \mathcal{Y}^n \rightarrow \mathcal{P}(\mathcal{X}^n)$. A (n, R, Δ, ϵ) -code, minimal error probability and error exponent are defined in a standard manner. The latter two are denoted by $\varepsilon_{X|Y}(n, R, \Delta)$ and $E_{X|Y}(R, \Delta)$ respectively.

This problem is very similar to its counterpart with no side information, i.e. given $\mathbf{Y} = \mathbf{y}$, the problem reduces to encoding and decoding a memoryless source (not necessarily i.i.d.) with distribution $P_{X|Y}(\mathbf{x}|\mathbf{y}) = \prod_{i=1}^n P_{X|Y}(x_i|y_i)$. Nevertheless, it is still useful to characterize the error exponent in this case, as it provides an upper bound for the more interesting case with decoder side information only. To this end, we define $F_{X|Y}(R)$ on $0 \leq R < \log |\mathcal{S}(P_X)|$ as

$$F_{X|Y}(R) \triangleq \min_{Q_{XY}: H(Q_{X|Y}|Q_Y) \geq R} D(Q_{XY} \| P_{XY}). \quad (21)$$

$F_{X|Y}(R)$ is continuous, convex and increasing on its domain; and is zero for $0 \leq R \leq H(P_{X|Y}|P_Y)$. Moreover, we have $F_{X|Y}(R) = E_{X|Y}(R, 0)$, that is the error exponent in the almost lossless case. $F_{X|Y}(R)$ also admits the following dual form in terms of Arimoto's conditional Rényi entropy [11]

$$F_{X|Y}(R) = \sup_{\rho \geq 0} \rho \left(R - H_{\frac{1}{1+\rho}}(X|Y) \right) \quad (22)$$

obtained using Lagrangian duality techniques (see, e.g., the proof of [10, Equation (32)] by Bunte and Lapidoth).

Theorem 2. *Let (R, Δ) be a rate-distortion pair such that $H(P_{X|Y}|P_Y) < R + \Delta < \log |\mathcal{S}(P_X)|$. Then*

$$E_{X|Y}(R, \Delta) = F_{X|Y}(R + \Delta). \quad (23)$$

The proof of Theorem 2 (omitted for brevity) is very similar to that of Theorem 1, but relies on conditional types.

¹By continuity and since $\bigcup_{n \in \mathbb{N}} \mathcal{P}_n(\mathcal{X})$ is dense in $\mathcal{P}(\mathcal{X})$ [9, Rem. 2].

B. Decoder side information: Wyner-Ziv

We now turn our attention to the case where the side information sequence \mathbf{Y} is only available at the decoder. This is the Wyner-Ziv setting, specialized to the log-loss distortion measure. A lossy source code of block-length n here is given by the pair $\phi_n : \mathcal{X}^n \rightarrow \mathcal{M}_n$ and $\varphi_n : \mathcal{M}_n \times \mathcal{Y}^n \rightarrow \mathcal{P}(\mathcal{X}^n)$. The minimal error probability and error exponent are denoted by $\varepsilon_{X|Y}^{\text{WZ}}(n, R, \Delta)$ and $E_{X|Y}^{\text{WZ}}(R, \Delta)$ respectively.

Next, we observe that the encoder-decoder side information result in Theorem 2 provides the following converse bound

$$E_{X|Y}^{\text{WZ}}(R, \Delta) \leq E_{X|Y}^{\text{SP}}(R, \Delta) \triangleq F_{X|Y}(R + \Delta). \quad (24)$$

For $\Delta = 0$, the setting reduces to the Slepian-Wolf problem, and we denote the error exponent by $E_{X|Y}^{\text{SW}}(R)$. The bound $E_{X|Y}^{\text{SW}}(R) \leq F_{X|Y}(R)$, a special case of (24), was obtained by Gallager in [12] (see also Csiszár and Körner [13, Theorem 3]). This bound on the Slepian-Wolf error exponent is sometimes referred to as the sphere-packing exponent [14], due to close resemblance to the sphere-packing exponent in channel coding. Similarly, the bound in (24) can be thought of as a sphere-packing exponent for the log-loss Wyner-Ziv setting.

Next, we derive a lower bound for $E_{X|Y}^{\text{WZ}}(R, \Delta)$. Define the function $\tilde{F}_{X|Y}(R)$ on $0 \leq R < \log |\mathcal{S}(P_X)|$ as

$$\tilde{F}_{X|Y}(R) \triangleq \min_{Q_{XY}} \left\{ D(Q_{XY} \| P_{XY}) + |R - H(Q_{X|Y}|Q_Y)|^+ \right\}$$

where $|a|^+ \triangleq \max\{0, a\}$. Note that $\tilde{F}_{X|Y}(R) \leq F_{X|Y}(R)$. Moreover, $\tilde{F}_{X|Y}(R)$ admits a dual form

$$\tilde{F}_{X|Y}(R) = \max_{\rho \in [0, 1]} \rho \left(R - H_{\frac{1}{1+\rho}}(X|Y) \right). \quad (25)$$

$\tilde{F}_{X|Y}(R)$ is an achievable error exponent in the Slepian-Wolf setting [12], [13], referred to as the random-coding error exponent, as it is achieved through random coding (or binning) in close resemblance to the random-coding exponent in channel coding. A corresponding random-coding error exponent for the log-loss Wyner-Ziv problem is presented next.

Theorem 3. *Let (R, Δ) be a rate-distortion pair such that $H(P_{X|Y}|P_Y) < R + \Delta < \log |\mathcal{S}(P_X)|$. Then*

$$E_{X|Y}^{\text{WZ}}(R, \Delta) \geq E_{X|Y}^{\text{r}}(R, \Delta) \triangleq \tilde{F}_{X|Y}(R + \Delta).$$

For fixed Δ , the exponents $E_{X|Y}^{\text{r}}(R, \Delta)$ and $E_{X|Y}^{\text{SP}}(R, \Delta)$ coincide on $H(P_{X|Y}|P_Y) - \Delta < R \leq R_{\text{cr}}$, where R_{cr} is the largest rate at which the convex curve $E_{X|Y}^{\text{SP}}(R, \Delta)$, as a function of R , meets its supporting line of slope 1. Note that R_{cr} is reminiscent of the *critical rate* in channel coding. Above this rate, the two exponents differ in general.

C. Proof of Theorem 3

The proof is based on random binning and a list decoding variant of the universal minimum entropy decoding rule [13], [15]. In the analysis of this scheme, we use $H(\mathbf{x}|\mathbf{y})$ as a shorthand notation for the conditional entropy $H(P_{\mathbf{x}|\mathbf{y}}|P_{\mathbf{y}})$ calculated from the joint type $P_{\mathbf{x}\mathbf{y}} = P_{\mathbf{x}|\mathbf{y}}P_{\mathbf{y}}$.

Let $J_n = |\mathcal{P}_n(\mathcal{X})|$ and $M_n = \lfloor (1+n)^{-|\mathcal{X}|} e^{nR} \rfloor$ for fixed R and n . A binning function $b_n : \mathcal{X}^n \rightarrow \{1, 2, \dots, M_n\}$ is a mapping that assigns an index $b_n(\mathbf{x}) \in \{1, 2, \dots, M_n\}$ to every source sequence $\mathbf{x} \in \mathcal{X}^n$. For a fixed bin assignment, determined by a given binning function b_n , the set of all source sequences with the same bin index as \mathbf{x} is denoted by

$$\mathcal{B}_n(\mathbf{x}|b_n) \triangleq \{\hat{\mathbf{x}} \in \mathcal{X}^n : b_n(\hat{\mathbf{x}}) = b_n(\mathbf{x})\}.$$

Further on, we will analyze the error probability averaged over an ensemble of binning functions. To that end, we denote a random binning function by B_n , where b_n is a realization of B_n . We use Gallager's ensemble [12]: every source sequence is assigned a bin index uniformly at random; and bin assignments are pairwise independent across sequences.

Encoding: The generated source sequence \mathbf{x} is encoded into $\phi_n(\mathbf{x}) = (t_n(\mathbf{x}), b_n(\mathbf{x}))$, where $t_n(\mathbf{x})$ is the type index and $b_n(\mathbf{x})$ is the bin index. The rate satisfies $\frac{1}{n} \log(J_n M_n) \leq R$.

Decoding: Upon receiving $\phi_n(\mathbf{x})$, the decoder knows that \mathbf{x} is in the set $\mathcal{T}_n(P_{\mathbf{x}}) \cap \mathcal{B}_n(\mathbf{x}|b_n)$. Now suppose that the side information sequence is equal to \mathbf{y} . For every sequence $\hat{\mathbf{x}} \in \mathcal{T}_n(P_{\mathbf{x}}) \cap \mathcal{B}_n(\mathbf{x}|b_n)$, the decoder computes the conditional entropy $H(\hat{\mathbf{x}}|\mathbf{y})$ and produces a list $\mathcal{L}_n(\phi_n(\mathbf{x}), \mathbf{y})$ of size

$$|\mathcal{L}_n(\phi_n(\mathbf{x}), \mathbf{y})| = \min \{ \lfloor e^{n\Delta} \rfloor, |\mathcal{T}_n(P_{\mathbf{x}}) \cap \mathcal{B}_n(\mathbf{x}|b_n)| \}$$

comprising source sequences with the lowest conditional entropy. The soft reconstruction is taken to be uniformly supported on $\mathcal{L}_n(\phi_n(\mathbf{x}), \mathbf{y})$. It is clear that an excess distortion error occurs if $\mathcal{L}_n(\phi_n(\mathbf{x}), \mathbf{y})$ does not include the encoded source sequence \mathbf{x} . Moreover, by setting $\Delta = 0$, we recover the classical minimum entropy decoder.

Error probability: Let $\mathcal{E}_n(\mathbf{x}|\mathbf{y}, b_n)$ denote the set of all source sequences other than \mathbf{x} , but with the same type and bin as \mathbf{x} , and a conditional entropy smaller than or equal to that of \mathbf{x} given \mathbf{y} . For an excess distortion error to occur, we must have $|\mathcal{E}_n(\mathbf{x}|\mathbf{y}, b_n)| \geq \lfloor e^{n\Delta} \rfloor = e^{n(\Delta - \delta_n)}$, where $\delta_n \geq 0$ goes to zero as n grows large. The excess distortion error probability, averaged over B_n , is hence bounded above by

$$\mathbb{P} \left[|\mathcal{E}_n(\mathbf{X}|\mathbf{Y}, B_n)| \geq e^{n(\Delta - \delta_n)} \right] \leq \sum_{\mathbf{x}, \mathbf{y}} P_{XY}^n(\mathbf{x}, \mathbf{y}) \min \left\{ 1, e^{-n(\Delta - \delta_n)} \mathbb{E} [|\mathcal{E}_n(\mathbf{x}|\mathbf{y}, B_n)|] \right\} \quad (26)$$

which follows from Markov's inequality combined with the trivial upper bound of 1. We now take a small detour to bound $\mathbb{E} [|\mathcal{E}_n(\mathbf{x}|\mathbf{y}, B_n)|]$. To this end, define the set $\mathcal{E}'_n(\mathbf{x}|\mathbf{y}) \triangleq \{\hat{\mathbf{x}} \in \mathcal{T}_n(P_{\mathbf{x}}) : \hat{\mathbf{x}} \neq \mathbf{x}, H(\hat{\mathbf{x}}|\mathbf{y}) \leq H(\mathbf{x}|\mathbf{y})\}$, and observe that

$$\mathbb{E} [|\mathcal{E}_n(\mathbf{x}|\mathbf{y}, B_n)|] = \sum_{\hat{\mathbf{x}} \in \mathcal{E}'_n(\mathbf{x}|\mathbf{y})} \mathbb{P} [B_n(\hat{\mathbf{x}}) = B_n(\mathbf{x})] = \frac{|\mathcal{E}'_n(\mathbf{x}|\mathbf{y})|}{M_n}$$

which follows from uniform pairwise independent bin assignment. Next, we note that

$$\begin{aligned} |\mathcal{E}'_n(\mathbf{x}|\mathbf{y})| &\leq \sum_{Q_{X|Y} \in \mathcal{P}_n(\mathcal{X}|\mathcal{Y}) : H(Q_{X|Y}|P_{\mathbf{y}}) \leq H(\mathbf{x}|\mathbf{y})} |\mathcal{T}_n(Q_{X|Y}|\mathbf{y})| \\ &\leq (n+1)^{|\mathcal{X}||\mathcal{Y}|} e^{nH(\mathbf{x}|\mathbf{y})}. \end{aligned} \quad (27)$$

From the above and $M_n = \lfloor (1+n)^{-|\mathcal{X}|} e^{nR} \rfloor$, we obtain

$$\mathbb{E} [|\mathcal{E}_n(\mathbf{x}|\mathbf{y}, B_n)|] \leq e^{nH(\mathbf{x}|\mathbf{y}) - nR + n\delta'_n} \quad (28)$$

from some $\delta'_n \geq 0$ which goes to zero as n grows large. Defining $\delta''_n \triangleq \delta_n + \delta'_n$, and using (28), it follows that

$$\min \left\{ 1, e^{-n(\Delta - \delta_n)} \mathbb{E} [|\mathcal{E}_n(\mathbf{x}|\mathbf{y}, B_n)|] \right\} \leq e^{-n|R + \Delta - H(\mathbf{x}|\mathbf{y})| + n\delta''_n}.$$

Plugging this back into (26), and invoking the usual random coding argument of the existence of a code, we obtain

$$\begin{aligned} \varepsilon_{X|Y}^{\text{WZ}}(n, R, \Delta) &\leq \sum_{\mathbf{y} \in \mathcal{Y}^n} \sum_{\mathbf{x} \in \mathcal{X}^n} P_{XY}^n(\mathbf{x}, \mathbf{y}) e^{-n|R + \Delta - H(\mathbf{x}|\mathbf{y})| + n\delta''_n} \\ &\leq \sum_{Q_{X|Y} \in \mathcal{P}_n(\mathcal{X}|\mathcal{Y})} e^{-nD(Q_{X|Y} \| P_{XY})} e^{-n|R + \Delta - H(Q_{X|Y}|Q_Y)| + n\delta''_n} \\ &\leq (n+1)^{|\mathcal{X}||\mathcal{Y}|} e^{-n\bar{F}_{X|Y}(R + \Delta) + n\delta''_n}. \end{aligned} \quad (29)$$

The result in Theorem 3 follows.

V. CONCLUDING REMARKS

In the high rate regime where $E_{X|Y}^r(R, \Delta)$ and $E_{X|Y}^{\text{sp}}(R, \Delta)$ diverge, it is possible to derive a tighter achievable exponent, which is a log-loss counterpart of the expurgated exponent in source coding [15]. We skip this due to lack of space. As an extension, it is of interest to derive error exponents and universal schemes for the multi-terminal settings in [1], [2].

REFERENCES

- [1] T. A. Courtade and R. D. Wesel, "Multiterminal source coding with an entropy-based distortion measure," in *Proc. IEEE Int. Symp. Inf. Theory*, 2011, pp. 2040–2044.
- [2] T. A. Courtade and T. Weissman, "Multiterminal source coding under logarithmic loss," *IEEE Trans. Inf. Theory*, vol. 60, no. 1, pp. 740–761, 2014.
- [3] Y. Y. Shkel and S. Verdú, "A single-shot approach to lossy source coding under logarithmic loss," *IEEE Trans. Inf. Theory*, vol. 64, no. 1, pp. 129–147, 2018.
- [4] Y. Shkel, M. Raginsky, and S. Verdú, "Universal lossy compression under logarithmic loss," in *Proc. IEEE Int. Symp. Inf. Theory*, 2017, pp. 1157–1161.
- [5] Y. Shkel, M. Raginsky, and S. Verdú, "Universal compression, list decoding, and logarithmic loss," in *Proc. IEEE Int. Symp. Inf. Theory*, 2018, pp. 206–210.
- [6] I. Csiszár and J. Körner, *Information theory: Coding theorems for discrete memoryless systems*. Cambridge University Press, 2011.
- [7] K. Marton, "Error exponent for source coding with a fidelity criterion," *IEEE Trans. Inf. Theory*, vol. 20, no. 2, pp. 197–199, 1974.
- [8] F. Jelinek, *Probabilistic information theory: Discrete and memoryless models*. McGraw-Hill, 1968.
- [9] G. Longo and A. Sgarro, "The source coding theorem revisited: A combinatorial approach," *IEEE Trans. Inf. Theory*, vol. 25, no. 5, pp. 544–548, 1979.
- [10] C. Bunte and A. Lapidoth, "Source coding, lists, and Rényi entropy," in *Proc. IEEE Inf. Theory Workshop*, 2013, pp. 1–5.
- [11] S. Arimoto, "Information measures and capacity of order α for discrete memoryless channels," in *Topics in Information Theory*, I. Csiszár and P. Elias, Eds. Amsterdam: North-Holland Publ. Co, 1977, pp. 41–52.
- [12] R. G. Gallager, "Source coding with side information and universal coding," M.I.T. LIDS, Tech. Rep., 1976.
- [13] I. Csiszár and J. Körner, "Towards a general theory of source networks," *IEEE Trans. Inf. Theory*, vol. 26, no. 2, pp. 155–165, 1980.
- [14] B. G. Kelly and A. B. Wagner, "Improved source coding exponents via Witsenhausen's rate," *IEEE Trans. Inf. Theory*, vol. 57, no. 9, pp. 5615–5633, 2011.
- [15] I. Csiszár and J. Körner, "Graph decomposition: A new key to coding theorems," *IEEE Trans. Inf. Theory*, vol. 27, no. 1, pp. 5–12, 1981.