

Randomized Quantization and Source Coding with Constrained Output Distribution

Naci Saldi, Tamás Linder, Serdar Yüksel

Abstract—This paper studies fixed-rate randomized vector quantization under the constraint that the quantizer’s output has a given fixed probability distribution. A general representation of randomized quantizers that includes the common models in the literature is introduced via appropriate mixtures of joint probability measures on the product of the source and reproduction alphabets. Using this representation and results from optimal transport theory, the existence of an optimal (minimum distortion) randomized quantizer having a given output distribution is shown under various conditions. For sources with densities and the mean square distortion measure, it is shown that this optimum can be attained by randomizing quantizers having convex codecells. For stationary and memoryless source and output distributions a rate-distortion theorem is proved, providing a single-letter expression for the optimum distortion in the limit of large block-lengths.

Index Terms—Source coding, quantization, randomization, random coding, output-constrained distortion-rate function.

I. INTRODUCTION

A quantizer maps a source (input) alphabet into a finite collection of points (output levels) from a reproduction alphabet. A quantizer’s performance is usually characterized by its rate, defined as the logarithm of the number of output levels, and its expected distortion when the input is a random variable. One usually talks about randomized quantization when the quantizer used to encode the input signal is randomly selected from a given collection of quantizers. Although introducing randomization in the quantization procedure does not improve the optimal rate-distortion tradeoff, randomized quantizers may have other advantages over their deterministic counterparts.

In what appears to be the first work explicitly dealing with randomized quantization, Roberts [1] found that adding random noise to an image before quantization and subtracting the noise before reconstruction may result in a perceptually more pleasing image. Schuchman [2] and Gray and Stockham [3] analyzed versions of such so called *dithered* scalar quantizers where random noise (dither) is added to the input signal prior to uniform quantization. If the dither is subtracted after the quantization operation, the procedure is called subtractive dithering; otherwise it is called non-subtractive dithering.

The authors are with the Department of Mathematics and Statistics, Queen’s University, Kingston, ON, Canada, Email: {nsaldi,linder,yuksel}@mast.queensu.ca

This research was supported in part by the Natural Sciences and Engineering Research Council (NSERC) of Canada.

The material in this paper was presented in part at the 2013 IEEE International Symposium on Information Theory, Istanbul, Turkey, July 2013.

©2014 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

Under certain conditions, dithering results in uniformly distributed quantization noise that is independent of the input [2], [3], which allows a simple modeling of the quantization process by an additive noise channel. In the information theoretic literature the properties of entropy coded dithered lattice quantizers have been extensively studied. For example, such quantizers have been used to provide achievable bounds on the performance of universal lossy compression systems by Ziv [4] and Zamir and Feder [5], [6]. Recently Akyol and Rose [7], [8], introduced a class of randomized *nonuniform* scalar quantizers obtained via applying companding to a dithered uniform quantizer and investigated optimality conditions for the design of such quantizers. One should also note that the random codes used to prove the achievability part of Shannon’s rate-distortion theorem [9] can also be viewed as randomized quantizers.

Dithered uniform/lattice and companding quantizers, as well as random rate-distortion codes, pick a random quantizer from a “small” structured subset of all possible quantizers. Such special randomized quantizers may be suboptimal for certain tasks and one would like to be able to work with more general (or completely general) classes of randomized quantizers. For example, Li *et al.* [10] and Klejsa *et al.* [12] considered *distribution-preserving* dithered scalar quantization, where the quantizer output is restricted to have the same distribution as the source, to improve the perceptual quality of mean square optimal quantizers in audio and video coding. Dithered quantizers or other structured randomized quantizers classes likely do not provide optimal performance in this problem. In an unpublished work [11] the same authors considered more general distribution-preserving randomized vector quantizers and lower bounded the minimum distortion achievable by such schemes when the source is stationary and memoryless.

In this paper we propose a general model which formalizes the notion of randomly picking a quantizer from the set of *all* quantizers with a given number of output levels. Note that this set is much more complex and less structured than say the *parametric* family of all quantizers having a given number of convex codecells. Inspired by work in stochastic control (e.g., [13]) our model represents the set of all quantizers acting on a given source as a subset of all joint probability measures on the product of the source and reproduction alphabets. Then a randomized quantizer corresponds to a mixture of probability measures in this subset. The usefulness of the model is demonstrated by rigorously setting up a generalization of the distribution-preserving quantization problem where then the goal is to find a randomized quantizer minimizing the distortion under the constraint that the output has a given distribution (not necessarily that of the source). We show that under

quite general conditions an optimal solution (i.e., an optimal randomized quantizer) exists for this generalized problem. We also consider a relaxed version of the output distribution constraint where the output distribution is only required to belong to some neighborhood (in the weak topology) of a target distribution. For this problem we show the optimality of randomizing among finitely many quantizers. While for fixed quantizer dimension we can only provide existence results, for stationary and memoryless source and output distributions we also develop a rate-distortion theorem which identifies the minimum distortion in the limit of large block lengths in terms of the so-called output-constrained distortion-rate function. This last result solves a general version of a problem that was left open in [11].

The rest of the paper is organized as follows. In Section II we introduce our general model for randomized quantization and show its equivalence to other models more common in the information theoretic literature. In Section III the randomized quantization problem with an output distribution constraint is formulated and the existence of an optimal solution is shown using optimal transport theory. For the special but important case of sources with densities and the mean square distortion measure, we show that this optimum can be attained by randomizing quantizers having convex codecells. In Section IV a relaxed version of output distribution constraint is considered where finitely randomized quantizers are optimal. In Section V we present and prove a rate-distortion theorem for fixed-rate lossy source coding with an output distribution constraint. Many of the proofs are quite technical and they are relegated to the Appendix.

II. MODELS OF RANDOMIZED QUANTIZATION

A. Notation

In this paper X denotes the input alphabet and Y is the reconstruction (output) alphabet. Throughout the paper we set $X = Y = \mathbb{R}^n$, the n -dimensional Euclidean space for some $n \geq 1$, although most of the results hold in more general settings; for example if the input and output alphabets are Polish (complete and separable metric) spaces. If E is a metric space, $\mathcal{B}(E)$ and $\mathcal{P}(E)$ will denote the Borel σ -algebra on E and the set of probability measures on $(E, \mathcal{B}(E))$, respectively. It will be tacitly assumed that any metric space is equipped with its Borel σ -algebra and all probability measures on such spaces will be Borel measures. The product of metric spaces will be equipped with the product Borel σ -algebra. Unless otherwise specified, the term “measurable” will refer to Borel measurability. We always equip $\mathcal{P}(E)$ with the Borel σ -algebra $\mathcal{B}(\mathcal{P}(E))$ generated by the topology of weak convergence [14].

B. Three models of randomized quantization

An M -level quantizer (M is a positive integer) from the input alphabet X to the reconstruction alphabet Y is a Borel measurable function $q : X \rightarrow Y$ whose range $q(X) = \{q(x) : x \in X\}$ contains at most M points of Y . If \mathcal{Q}_M denotes the set of all M -level quantizers, then our definition implies $\mathcal{Q}_M \subset \mathcal{Q}_{M+1}$ for all $M \geq 1$.

In what follows we define three models of randomized quantization; two that are commonly used in the source coding literature, and a third abstract model that will nevertheless prove very useful.

Model 1

One general model of M -level randomized quantization that is often used in the information theoretic literature is depicted in Fig. 1.

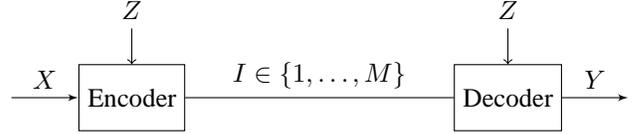


Figure 1. Randomized source code (quantizer)

Here X and Y are the source and output random variables taking values in X and Y , respectively. The index I takes values in $\{1, \dots, M\}$, and Z is a $Z = \mathbb{R}^m$ -valued random variable which is independent of X and which is assumed to be available at both the encoder and the decoder. The encoder is a measurable function $e : X \times Z \rightarrow \{1, \dots, M\}$ which maps (X, Z) to I , and the decoder is a measurable function $d : \{1, \dots, M\} \times Z \rightarrow Y$ which maps (I, Z) to Y . For a given source distribution, in a probabilistic sense a Model 1 quantizer is determined by the triple (e, d, ν) , where ν denotes the distribution of Z .

Note that codes used in the random coding proof of the forward part of Shannon’s rate distortion theorem can be realized as Model 1 quantizers. In this case Z may be taken to be the random codebook consisting of $M = 2^{nR}$ code vectors of dimension n , each drawn independently from a given distribution. This Z can be represented as an $m = nM$ -dimensional random vector that is independent of X . The encoder outputs the index I of the code vector Y in the codebook that best matches X (in distortion or in a joint-typicality sense) and the decoder can reconstruct this Y since it is a function of I and Z .

Model 2

Model 1 can be collapsed into a more tractable equivalent model. In this model, a randomized quantizer is a pair (q, ν) , where $q : X \times Z \rightarrow Y$ is a measurable mapping such that $q(\cdot, z)$ is an M -level quantizer for all $z \in Z$ and ν is a probability measure on Z , the distribution of the randomizing random variable Z . Here q is the composition of the encoder and the decoder in Model 1: $q(x, z) = d(e(x, z), z)$.

Model 2 quantizers include, as special cases, subtractive and non-subtractive dithering of M -level uniform quantizers, as well as the dithering of non-uniform quantizers. For example, if $n = m = 1$ and q_u denotes a uniform quantizer, then

$$q(x, z) = q_u(x + z) - z$$

is a dithered uniform quantizer using subtractive dithering,

$$q(x, z) = q_u(x + z)$$

is a dithered uniform quantizer with non-subtractive dithering, and with appropriate mappings g and h ,

$$q(x, z) = h(q_u(g(x) + z) - z).$$

is a dithered non-uniform quantizer (e.g., [10] and [8]). We note that dithered lattice quantizers [4], [5], [15] can also be considered as Model 2 type randomized quantizers when the source has a bounded support (so that with probability one only finitely many lattice points can occur as output points).

Let $\rho : X \times Y \rightarrow \mathbb{R}$ be a nonnegative measurable function, called the distortion measure. From now on we assume that the source X has distribution μ (denoted as $X \sim \mu$). The distortion associated with Model 2 quantizer (q, ν) or with Model 1 quantizer (e, d, ν) , with $q(x, z) = d(e(x, z), z)$, is the expectation

$$\begin{aligned} L(q, \nu) &= \int_Z \int_X \rho(x, q(x, z)) \mu(dx) \nu(dz) \\ &= E[\rho(X, q(X, Z))] \end{aligned} \quad (1)$$

where $Z \sim \nu$ is independent of X .

Model 3

In this model, instead of considering quantizers as functions that map X into a finite subset of Y , first we represent them as special probability measures on $X \times Y$ (see, e.g., [16], [17], [18], [19]). This leads to an alternative representation where a randomized quantizer is identified with a mixture of probability measures. In certain situations the space of these “mixing probabilities” representing *all* randomized quantizers will turn out to be more tractable than considering the quite unstructured space of all Model 1 triples (e, d, ν) or Model 2 pairs (q, ν) .

A *stochastic kernel* [20] (or *regular conditional probability* [21]) on Y given X is a function $Q(dy|x)$ such that for each $x \in X$, $Q(\cdot|x)$ is a probability measure on Y , and for each Borel set $B \subset Y$, $Q(B|\cdot)$ is a measurable function from X to $[0, 1]$. A quantizer q from X into Y can be represented as a stochastic kernel Q on Y given X by letting [17], [16],

$$Q(dy|x) = \delta_{q(x)}(dy),$$

where δ_u denotes the point mass at u : $\delta_u(A) = 1$ if $u \in A$ and $\delta_u(A) = 0$ if $u \notin A$ for any Borel set $A \subset Y$.

If we fix the distribution μ of the source X , we can also represent q by the probability measure $v(dx dy) = \mu(dx) \delta_{q(x)}(dy)$ on $X \times Y$. Thus we can identify the set \mathcal{Q}_M of all M -level quantizers from X to Y with the following subset of $\mathcal{P}(X \times Y)$:

$$\begin{aligned} \Gamma_\mu(M) & \\ &= \{v \in \mathcal{P}(X \times Y) : v(dx dy) = \mu(dx) \delta_{q(x)}(dy), q \in \mathcal{Q}_M\}. \end{aligned} \quad (2)$$

Note that $q \mapsto \mu(dx) \delta_{q(x)}(dy)$ maps \mathcal{Q}_M onto $\Gamma_\mu(M)$, but this mapping is one-to-one only if we consider equivalence classes of quantizers in \mathcal{Q}_M that are equal μ almost everywhere (μ -a.e.).

We equip $\mathcal{P}(X \times Y)$ with the topology of weak convergence (weak topology) which is metrizable with the Prokhorov metric, making $\mathcal{P}(X \times Y)$ into a Polish space [14]. The following lemma is proved in the Appendix A.

Lemma 1. $\Gamma_\mu(M)$ is a Borel subset of $\mathcal{P}(X \times Y)$.

Now we are ready to introduce Model 3 for randomized quantization. Let P be a probability measure on $\mathcal{P}(X \times Y)$ which is supported on $\Gamma_\mu(M)$, i.e., $P(\Gamma_\mu(M)) = 1$. Then P induces a “randomized quantizer” $v_P \in \mathcal{P}(X \times Y)$ via

$$v_P(A \times B) = \int_{\Gamma_\mu(M)} v(A \times B) P(dv)$$

for Borel sets $A \subset X$ and $B \subset Y$, which we abbreviate to

$$v_P = \int_{\Gamma_\mu(M)} v P(dv). \quad (3)$$

Since each v in $\Gamma_\mu(M)$ corresponds to a quantizer with input distribution μ , P can be thought as a probability measure on the set of all M -level quantizers \mathcal{Q}_M .

Let $\mathcal{P}_0(\Gamma_\mu(M))$ denote the set of probability measures on $\mathcal{P}(X \times Y)$ supported on $\Gamma_\mu(M)$. We define the set of M -level Model 3 randomized quantizers as

$$\begin{aligned} \Gamma_\mu^R(M) & \\ &= \left\{ v_P \in \mathcal{P}(X \times Y) : v_P = \int_{\Gamma_\mu(M)} v P(dv), P \in \mathcal{P}_0(\Gamma_\mu(M)) \right\}. \end{aligned} \quad (4)$$

Note that if $v_P \in \Gamma_\mu^R(M)$ is a Model 3 quantizer, then the X -marginal of v_P is equal to μ , and if X and Y are random vectors (defined on the same probability space) with joint distribution v_P , then they provide a stochastic representation of the random quantizer’s input and output, respectively. Furthermore, the distortion associated with v_P is

$$\begin{aligned} L(v_P) &:= \int_{X \times Y} \rho(x, y) v_P(dx dy) \\ &= \int_{\Gamma_\mu(M)} \int_{X \times Y} \rho(x, y) v(dx dy) P(dv) \\ &= E[\rho(X, Y)]. \end{aligned}$$

C. Equivalence of models

Here we show that the three models of randomized quantization are essentially equivalent. As before, we assume that the source distribution μ is fixed. The following two results are proved in Appendix B and Appendix C, respectively.

Theorem 1. For each Model 2 randomized quantizer (q, ν) there exists a Model 3 randomized quantizer $v_P \in \Gamma_\mu^R(M)$ such that $(X, Y) = (X, q(X, Z))$ has distribution v_P . Conversely, for any $v_P \in \Gamma_\mu^R(M)$ there exists a Model 2 randomized quantizer (q, ν) such that $(X, q(X, Z)) \sim v_P$.

Theorem 2. Models 1 and 2 of randomized quantization are equivalent in the sense of Theorem 1.

Remark 1.

- (a) Clearly, any two equivalent randomized quantizers have the same distortion. The main result of this section is Theorem 1. Theorem 2 is intuitively obvious, but proving that any Model 2 quantizer can be decomposed into an equivalent Model 1 quantizer with measurable encoder and decoder is not quite trivial.

- (b) Since the dimension m of the randomizing random vector Z was arbitrary, we can take $m = 1$ in Theorem 1. In fact, the proof also implies that any Model 2 or 3 randomized quantizer is equivalent (in the sense of Theorem 1) to a Model 2 quantizer (q, ν) , where $q : X \times [0, 1] \rightarrow Y$ and ν is the uniform distribution on $[0, 1]$.
- (c) Assume that (Z, \mathcal{A}, ν) is an *arbitrary* probability space. For any randomized quantizer $q : X \times Z \rightarrow Y$ in the form $q(X, Z)$, where $Z \sim \nu$ is independent of X , there exists a Model 3 randomized quantizer v_P such that $(X, q(X, Z)) \sim v_P$. This can be proved by using the same proof method as in Theorem 1. In view of the previous remark and Theorem 1, this means that uniform randomization over the unit interval $[0, 1]$ suffices under the most general circumstances.
- (d) All results in this section remain valid if the input and reproduction alphabets X and Y are arbitrary uncountable Polish spaces. In this case, uniform randomization over the unit interval still provides the most general model possible.

In the next two sections, Model 3 will be used to represent randomized quantizers because it is particularly suited to treating the optimal randomized quantization problem under an output distribution constraint.

III. OPTIMAL RANDOMIZED QUANTIZATION WITH FIXED OUTPUT DISTRIBUTION

Let ψ be a probability measure on Y and let $\Lambda(M, \psi)$ denote the set of all M -level Model 2 randomized quantizers (q, ν) such that the output $q(X, Z)$ has distribution ψ . As before, we assume that $X \sim \mu$, $Z \sim \nu$, and Z and X are independent. We want to show the existence of a minimum-distortion randomized quantizer having output distribution ψ , i.e., the existence of $(q^*, \nu^*) \in \Lambda(M, \psi)$ such that

$$L(q^*, \nu^*) = \inf_{(q, \nu) \in \Lambda(M, \psi)} L(q, \nu).$$

If we set $\psi = \mu$, the above problem is reduced to showing the existence of a distribution-preserving randomized quantizer [10], [11] having minimum distortion.

The set of M -level randomized quantizers is a fairly general (nonparametric) set of functions and it seems difficult to investigate the existence of an optimum directly. On the other hand, Model 3 provides a tractable framework for establishing the existence of an optimal randomized quantizer under quite general conditions.

Let $\Gamma_{\mu, \psi}$ be the set of all joint distributions $v \in P(X \times Y)$ having X -marginal μ and Y -marginal ψ . Then

$$\Gamma_{\mu, \psi}^R(M) = \Gamma_{\mu}^R(M) \cap \Gamma_{\mu, \psi} \quad (5)$$

is the subset of Model 3 randomized quantizers which corresponds to the class of output-distribution-constrained Model 2 randomized quantizers $\Lambda(M, \psi)$.

For any $v \in P(X \times Y)$ let

$$L(v) = \int_{X \times Y} \rho(x, y) v(dx dy).$$

Using these definitions, finding optimal randomized quantizers with a given output distribution can be posed as finding $v \in \Gamma_{\mu, \psi}^R(M)$ which minimizes $L(v)$, i.e.,

$$\begin{aligned} \text{(P1) minimize } & L(v) \\ \text{subject to } & v \in \Gamma_{\mu, \psi}^R(M). \end{aligned}$$

We can prove the existence of the minimizer for **(P1)** under either of the following assumptions. Here $\|x\|$ denotes the Euclidean norm of $x \in \mathbb{R}^n$.

ASSUMPTION 1: $\rho(x, y)$ is continuous and $\psi(B) = 1$ for some compact subset B of Y .

ASSUMPTION 2: $\rho(x, y) = \|x - y\|^2$.

Theorem 3. *Suppose $\inf_{v \in \Gamma_{\mu, \psi}^R(M)} L(v) < \infty$. Then there exists a minimizer with finite cost for problem **(P1)** under either Assumption 1 or Assumption 2.*

The theorem is proved in Appendix D with the aid of optimal transport theory [22]. The optimal transport problem for marginals $\pi \in \mathcal{P}(X)$, $\lambda \in \mathcal{P}(Y)$ and cost function $c : X \times Y \rightarrow [0, \infty]$ is defined as

$$\begin{aligned} \text{minimize } & \int_{X \times Y} c(x, y) v(dx dy) \\ \text{subject to } & v \in \Gamma_{\pi, \lambda}. \end{aligned}$$

In the proof of Theorem 3 we set up a relaxed version of the optimization problem **(P1)**. We show that if the relaxed problem has a minimizer, then **(P1)** also has a minimizer, and then prove the existence of a minimizer for the relaxed problem using results from optimal transport theory.

Remark 2. Note that the product distribution $\mu \otimes \psi$ corresponds to a 1-level randomized quantizer (the equivalent Model 2 randomized quantizer is given by $q(x, z) = z$ and $Z \sim \nu$). Hence $\mu \otimes \psi \in \Gamma_{\mu, \psi}^R(M)$ for all $M \geq 1$, and if $L(\mu \otimes \psi) < \infty$, then the condition $\inf_{v \in \Gamma_{\mu, \psi}^R(M)} L(v) < \infty$ holds. In particular, if both μ and ψ have finite second moments $\int \|x\|^2 \mu(dx) < \infty$ and $\int \|y\|^2 \psi(dy) < \infty$, and $\rho(x, y) = \|x - y\|^2$ (Assumption 2), then $\inf_{v \in \Gamma_{\mu, \psi}^R(M)} L(v) < \infty$.

Optimal transport theory can also be used to show that, under some regularity conditions on the input distribution and the distortion measure, the randomization can be restricted to quantizers having a certain structure. Here we consider sources with densities and the mean square distortion. A quantizer $q : X \rightarrow Y$ with output points $q(X) = \{y_1, \dots, y_k\} \subset Y$ is said to have *convex codecells* if $q^{-1}(y_i) = \{x : q(x) = y_i\}$ is a convex subset of $X = \mathbb{R}^n$ for all $i = 1, \dots, k$. Let $\mathcal{Q}_{M, c}$ denote the set of all M -level quantizers having convex codecells. The proof of the following theorem is given in Appendix E.

Theorem 4. *Suppose $\rho(x, y) = \|x - y\|^2$ and μ admits a probability density function. Then an optimal randomized quantizer in Theorem 3 can be obtained by randomizing over quantizers with convex cells. That is*

$$\min_{v \in \Gamma_{\mu, \psi}^R(M)} L(v) = \min_{v \in \Gamma_{\mu, \psi}^{R, c}(M)} L(v),$$

where $\Gamma_{\mu, \psi}^{R, c}(M)$ represents the Model 3 quantizers with output distribution ψ that are obtained by replacing \mathcal{Q}_M with $\mathcal{Q}_{M, c}$ in (2).

Remark 3. Each quantizer having M convex codecells can be described using $nM + (n + 1)M(M - 1)/2$ real parameters if μ has a density and any two quantizers that are μ -a.e. equal are considered equivalent. One obtains such a parametric description by specifying the M output points using nM real parameters, and specifying the M convex polytopal codecells by $M(M - 1)/2$ hyperplanes separating pairs of distinct codecells using $(n + 1)M(M - 1)/2$ real parameters. Thus Theorem 4 replaces the *nonparametric* family of quantizers Q_M in Theorem 3 with the *parametric* family $Q_{M,c}$.

IV. APPROXIMATION WITH FINITE RANDOMIZATION

Since randomized quantizers require common randomness that must be shared between the encoder and the decoder, it is of interest to see how one can approximate the optimal cost by randomizing over finitely many quantizers. Clearly, if the target probability measure ψ on Y is not finitely supported, then no finite randomization exists with this output distribution. In this section we relax the fixed output distribution constraint and consider the problem where the output distribution belongs to some neighborhood (in the weak topology) of ψ . We show that one can always find a finitely randomized quantizer which is optimal (resp., ε -optimal) for this relaxed problem if the distortion measure is continuous and bounded (resp., arbitrary).

Let $B(\psi, \delta)$ denote the open ball in $\mathcal{P}(Y)$, with respect to the Prokhorov metric [14] (see also (20) in Appendix F), having radius $\delta > 0$ and centered at the target input distribution ψ . Also, let $\mathcal{M}_{\mu, \psi}^{\delta}$ denote the set of all $v \in \Gamma_{\mu}^{\text{R}}(M)$ whose Y marginal belongs to $B(\psi, \delta)$. That is, $\mathcal{M}_{\mu, \psi}^{\delta}$ represents all randomized quantizers in $\Gamma_{\mu}^{\text{R}}(M)$ whose output distribution is within distance δ of the target distribution ψ . We consider the following relaxed version of the minimization problem **(P1)**:

$$\begin{aligned} \text{(P3)} \quad & \text{minimize } L(v) \\ & \text{subject to } v \in \mathcal{M}_{\mu, \psi}^{\delta}. \end{aligned}$$

The set of *finitely randomized* quantizers in $\Gamma_{\mu}^{\text{R}}(M)$ is obtained by taking finite mixtures of quantizers in $\Gamma_{\mu}(M)$, i.e.,

$$\begin{aligned} & \Gamma_{\mu}^{\text{FR}}(M) \\ & = \left\{ v_P \in \Gamma_{\mu}^{\text{R}}(M) : v_P = \int_{\Gamma_{\mu}(M)} v P(dv), |\text{supp}(P)| < \infty \right\}. \end{aligned}$$

Theorem 5. *Assume the distortion measure ρ is continuous and bounded and let $v \in \mathcal{M}_{\mu, \psi}^{\delta}$ be arbitrary. Then there exists v_F in $\mathcal{M}_{\mu, \psi}^{\delta} \cap \Gamma_{\mu}^{\text{FR}}(M)$ such that $L(v_F) \leq L(v)$.*

The proof is given in Appendix F.

Although the minimum in **(P3)** may not be achieved by any $v \in \mathcal{M}_{\mu, \psi}^{\delta}$, the theorem implies that if the problem has a solution, it also has a solution in the set of finitely randomized quantizers.

Corollary 1. *Assume ρ is continuous and bounded and suppose there exists $v^* \in \mathcal{M}_{\mu, \psi}^{\delta}$ with $L(v^*) = \inf_{v \in \mathcal{M}_{\mu, \psi}^{\delta}} L(v)$. Then there exists $v_F \in \mathcal{M}_{\mu, \psi}^{\delta} \cap \Gamma_{\mu}^{\text{FR}}(M)$ such that $L(v_F) = L(v^*)$.*

The continuity of L , implied by the boundedness and continuity of ρ is crucial in the proof of Theorem 5 and thus for Corollary 1. However, the next theorem shows that for an arbitrary ρ , any $\varepsilon > 0$, and $v \in \mathcal{M}_{\mu, \psi}^{\delta}$, there exists v_F in $\mathcal{M}_{\mu, \psi}^{\delta} \cap \Gamma_{\mu}^{\text{FR}}(M)$ such that $L(v_F) \leq L(v) + \varepsilon$. That is, for any $\varepsilon > 0$ there exists an ε -optimal finitely randomized quantizer for **(P3)**. The theorem is proved in Appendix G

Theorem 6. *Let ρ be an arbitrary distortion measure and assume $\inf_{v \in \mathcal{M}_{\mu, \psi}^{\delta}} L(v) < \infty$. Then,*

$$\inf_{v \in \mathcal{M}_{\mu, \psi}^{\delta} \cap \Gamma_{\mu}^{\text{FR}}(M)} L(v) = \inf_{v \in \mathcal{M}_{\mu, \psi}^{\delta}} L(v).$$

Remark 4. The above results on finite randomization heavily depend on our use of the Prokhorov metric as a measure of “distance” between two probability measures. In particular, if one considers other measures of closeness, such as the Kullback-Leibler (KL) divergence or the total variation distance, then finite randomization may not suffice if the target output distribution is not discrete. In particular, if the target output distribution ψ has a density and $\tilde{\psi}$ denotes the (necessarily discrete) output distribution of any finitely randomized quantizer, then $\tilde{\psi}$ is not absolutely continuous with respect to ψ and for the KL divergence we have $D_{KL}(\tilde{\psi} \parallel \psi) = \infty$, while for the total variation distance we have $\|\tilde{\psi} - \psi\|_{TV} = 1$.

V. A SOURCE CODING THEOREM

After proving the existence of an optimum randomized quantizer in problem **(P1)** in Section IV, one would also like to evaluate the minimum distortion

$$L^* := \min\{L(v) : v \in \Gamma_{\mu, \psi}^{\text{R}}(M)\} \quad (6)$$

achievable for fixed source and output distributions μ and ψ and given number of quantization levels M . For any given blocklength n this seems to be a very hard problem in general. However, we are able to prove a rate-distortion type result that explicitly identifies L^* in the limit of large block lengths n if the source and output distributions correspond to two stationary and memoryless (i.e., i.i.d.) processes.

With a slight abuse of the notation used in previous sections, we let $X = Y = \mathbb{R}$ and consider a sequence of problems **(P1)** with input and output alphabets $X^n = Y^n = \mathbb{R}^n$, $n \geq 1$, and corresponding source and output distributions $\mu^n = \mu \otimes \dots \otimes \mu$ and $\psi^n = \psi \otimes \dots \otimes \psi$, the n -fold products of a two fixed probability measures $\mu, \psi \in \mathcal{P}(\mathbb{R})$. To emphasize the changing block length, $x^n = (x_1, \dots, x_n)$ and $y^n = (y_1, \dots, y_n)$ will denote generic elements of X^n and Y^n , respectively.

ASSUMPTION 3: The distortion measure is the average squared error given by

$$\rho_n(x^n, y^n) = \frac{1}{n} \sum_{i=1}^n \rho(x_i, y_i)$$

with $\rho(x, y) = (x - y)^2$. We assume that μ and ψ have finite second moments, i.e., $\int x^2 \mu(dx) < \infty$, $\int y^2 \psi(dy) < \infty$.

For $R \geq 0$ let $\Gamma_{\mu^n, \psi^n}^{\text{R}}(2^{nR})$ denote the set of n -dimensional Model 3 randomized quantizers defined in (5) having input

distribution μ^n , output distribution ψ^n , and at most 2^{nR} levels (i.e., rate R). Then

$$L_n(\mu, \psi, R) := \inf\{L(v) : v \in \Gamma_{\mu^n, \psi^n}^R(2^{nR})\}$$

is the minimum distortion achievable by such quantizers.

We also define

$$D(\mu, \psi, R) = \inf\{E[\rho(X, Y)] : X \sim \mu, Y \sim \psi, I(X; Y) \leq R\},$$

where the infimum is taken over pairs of all joint distributions of real random variables X and Y such that X has distribution μ , Y has distribution ψ , and their mutual information $I(X; Y)$ is upper bounded by R .

One can trivially adapt the standard proof from rate-distortion theory to show that similar to the distortion-rate function, $D(\mu, \psi, R)$ is a convex and nonincreasing function of R . Note that $D(\mu, \psi, R)$ is finite for all $R \geq 0$ by the assumption that μ and ψ have finite second moments. The distortion-rate function $D(\mu, R)$ of the i.i.d. source μ , is obtained from $D(\mu, \psi, R)$ as

$$D(\mu, R) = \inf_{\psi \in \mathcal{P}(Y)} D(\mu, \psi, R).$$

By a standard argument one can easily show that the sequence $\{nL_n(\mu, \psi, R)\}_{n \geq 1}$ is subadditive and so $\inf_{n \geq 1} L_n(\mu, \psi, R) = \lim_{n \rightarrow \infty} L_n(\mu, \psi, R)$. Thus the limit represents the minimum distortion achievable with rate- R randomized quantizers for an i.i.d. source with marginal μ under the constraint that the output is i.i.d. with marginal ψ . The next result proves that this limit is equal to $D(\mu, \psi, R)$, which one could thus call the ‘‘output-constrained distortion-rate function.’’

Theorem 7. *We have*

$$\lim_{n \rightarrow \infty} L_n(\mu, \psi, R) = D(\mu, \psi, R). \quad (7)$$

Remark 5.

- (a) As usual, the proof of the theorem consists of a converse and an achievability part. The converse (Lemma 2 below) directly follows from the usual proof of the converse part of the rate-distortion theorem. In fact, this was first noticed in [11] where the special case $\psi = \mu$ was considered and (in a different formulation) it was shown that for all n

$$L_n(\mu, \mu, R) \geq D(\mu, \mu, R).$$

Virtually the same argument implies that $L_n(\mu, \psi, R) \geq D(\mu, \psi, R)$ for all n and ψ . Nevertheless, we write out the proof in Appendix H since, strictly speaking, the proof in [11] is only valid if ψ is discrete with finite (Shannon) entropy or it has a density and finite differential entropy.

- (b) The proof of the converse part (Lemma 2) is valid for any randomized quantizer whose output Y^n satisfies $Y_i \sim \psi$, $i = 1, \dots, n$. Thus the theorem also holds if in the definition of $L_n(\mu, \psi, R)$, the randomized quantizers are required to have outputs with identically distributed (but not necessarily independent) components having common distribution ψ .

- (c) In [11] it was left as an open problem if $D(\mu, \mu, R)$ can be asymptotically achieved by a sequence of distribution-preserving randomized quantizers. The authors presented an incomplete achievability proof for the special case of Gaussian μ using dithered lattice quantization. We prove the achievability of $D(\mu, \psi, R)$ for arbitrary μ and ψ using a fundamentally different (but essentially non-constructive) approach. In particular, our proof is based on random coding where the codewords are uniformly distributed on the type class of an n -type that well approximates the target distribution ψ , combined with optimal coupling from mass transport theory.
- (d) With only minor changes in the proof, the theorem remains valid if $X = Y$ are arbitrary Polish spaces with metric d and $\rho(x, y) = d(x, y)^p$ for some $p \geq 1$. In this case the finite second moment conditions translate into $\int d(x, x_0)^p \mu(dx) < \infty$ and $\int d(y, y_0)^p \psi(dy) < \infty$ for some (and thus all) $x_0, y_0 \in X$.

Proof of Theorem 7. In this proof we use Model 2 of randomized quantization which is more suitable here than Model 3. Also, it is easier to deal with the rate-distortion performance than with the distortion-rate performance. Thus, following the notation in [23], for $D \geq 0$ we define the *minimum mutual information with constraint output* ψ as

$$I_m(\mu \| \psi, D) = \inf\{I(X; Y) : X \sim \mu, Y \sim \psi, E[\rho(X, Y)] \leq D\},$$

where the infimum is taken over pairs of all joint distributions of X with marginal μ and Y with marginal ψ such that $E[\rho(X, Y)] \leq D$. If this set of joint distributions is empty, we let $I_m(\mu \| \psi, D) = \infty$. Clearly, the extended real valued functions $I_m(\mu \| \psi, \cdot)$ and $D(R, \mu, \cdot)$ are inverses of each other. Hence $I_m(\mu \| \psi, D)$ is a nonincreasing, convex function of D .

The converse part of the theorem, i.e., the statement $L_n(\mu, \psi, R) \geq D(R, \mu, \psi)$ for all $n \geq 1$, is directly implied by the following lemma. The proofs of all lemmas in this section are given in Appendix H.

Lemma 2. *For all $n \geq 1$ if a randomized quantizer has input distribution μ^n , output distribution ψ^n , and distortion D , then its rate is lower bounded as*

$$R \geq I_m(\mu \| \psi, D).$$

In the rest of the proof we show the achievability of $D(R, \mu, \psi)$. We first prove this for finite alphabets and then generalize to continuous alphabets.

Let $X = Y$ be finite sets and assume that $\rho(x, y) = d(x, y)^p$, where d is a metric on X and $p > 0$. For each n let ψ_n be a closest n -type [24, Chapter 11] to ψ in the l_1 -distance which is absolutely continuous with respect to ψ , i.e., $\psi_n(y) = 0$ whenever $\psi(y) = 0$. Let D be such that $I_m(\mu \| \psi, D) < \infty$, let $\varepsilon > 0$ be arbitrary, and set $R = I_m(\mu \| \psi, D) + \varepsilon$. Assume $X^n \sim \mu^n$ for $n \geq 1$. For each n generate 2^{nR} codewords uniformly and independently drawn from $T_n(\psi_n)$, the *type class* of ψ_n [24], i.e., independently (of each other and of X^n) generate random codewords $U^n(1), \dots, U^n(2^{nR})$ such

that $U^n(i) \sim \psi_n^{(n)}$, where

$$\psi_n^{(n)}(y^n) = \begin{cases} \frac{1}{|T_n(\psi_n)|}, & \text{if } y^n \in T_n(\psi_n) \\ 0, & \text{otherwise.} \end{cases}$$

(As usual, for simplicity we assume that 2^{nR} is an integer.) Let \hat{X}^n denote the output of the nearest neighborhood encoder: $\hat{X}^n = \arg \min_{1 \leq i \leq 2^{nR}} \rho_n(X^n, U^n(i))$. In case of ties, we choose $U^n(i)$ with the smallest index i . The next lemma states the intuitively clear fact that \hat{X}^n is uniformly distributed on $T_n(\psi_n)$.

Lemma 3. $\hat{X}^n \sim \psi_n^{(n)}$.

The idea for this random coding scheme comes from [23] where an infinite i.i.d. codebook $\{U^n(i)\}_{i=1}^\infty$ was considered and the coding rate was defined as $(1/n) \log N_n$, where N_n is the smallest index i such that $\rho_n(X^n, U^n(i)) \leq D$. If the $U^n(i)$ are uniformly chosen from the type class $T_n(\psi_n)$, then by Theorem 1 and Appendix A and B of [23], $(1/n) \log N_n - I_m(\mu \| \psi_n, D) \rightarrow 0$ in probability.

Our scheme converts this variable-length random coding scheme into a fixed-rate scheme by considering, for each blocklength n , the finite codebook $\{U^n(i)\}_{i=1}^{2^{nR}}$. Letting $\rho_{\max} = \max_{x,y} \rho(x,y)$, the expected distortion of our scheme is bounded as

$$E[\rho_n(X^n, \hat{X}^n)] \leq D + \rho_{\max} \Pr\left\{\frac{1}{n} \log N_n > R\right\}.$$

Since $I_m(\mu \| \psi_n, D) \rightarrow I_m(\mu \| \psi, D)$ by the continuity of $I_m(\mu \| \psi, D)$ in ψ (see [23, Appendix A]), we have $R \geq I_m(\mu \| \psi_n, D) + \delta$ for some $\delta > 0$ if n is large enough. Thus the above bound implies

$$\limsup_{n \rightarrow \infty} E[\rho_n(X^n, \hat{X}^n)] \leq D. \quad (8)$$

Hence our random coding scheme has the desired rate and distortion as $n \rightarrow \infty$. However, its output \hat{X}^n has distribution $\psi_n^{(n)}$ instead of the required ψ^n . The next lemma shows that the normalized Kullback-Leibler divergence (relative entropy, [24]) between $\psi_n^{(n)}$ and ψ^n asymptotically vanishes.

Lemma 4. $\frac{1}{n} \mathcal{D}(\psi_n^{(n)} \| \psi^n) \rightarrow 0$ as $n \rightarrow \infty$.

Let $\pi, \lambda \in \mathcal{P}(\mathsf{X})$. The optimal transportation cost $\hat{T}_n(\pi, \lambda)$ between π and λ (see, e.g., [22]) with cost function ρ_n is defined by

$$\hat{T}_n(\pi, \lambda) = \inf\{E[\rho_n(U^n, V^n)] : U^n \sim \pi, V^n \sim \lambda\}, \quad (9)$$

where the infimum is taken over all joint distribution of pairs of random vectors (U^n, V^n) satisfying the given marginal distribution constraints. The joint distribution achieving $\hat{T}_n(\pi, \lambda)$ as well as the resulting pair (U^n, V^n) are both called an optimal coupling of π and λ . Optimal couplings exist when X is finite or $\mathsf{X} = \mathbb{R}^n$, $\rho(x, y) = (x - y)^2$, and both π and λ both have finite second moments [22].

Now consider an optimal coupling (\hat{X}^n, Y^n) of $\psi_n^{(n)}$ and ψ^n . If Z_1 and Z_2 are uniform random variables on $[0, 1]$ such that $Z = (Z_1, Z_2)$ is independent of X^n , then the random code and optimal coupling can be “realized” as

$(U^n(1), \dots, U^n(2^{nR})) = f_n(Z_1)$, $\hat{X}^n = \hat{f}_n(X^n, Z_1)$, and $Y^n = g_n(\hat{X}^n, Z_2)$, where f_n , \hat{f}_n , and g_n are suitable (measurable) functions. Combining random coding with optimal coupling this way gives rise to a randomized quantizer of type Model 2 whose output has the desired distribution ψ^n (see Fig. 2).

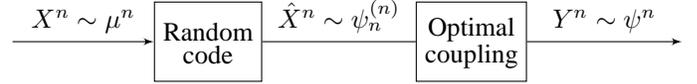


Figure 2. $D(R, \mu, \psi)$ achieving randomized quantizer scheme.

The next lemma uses Marton’s inequality [25] to show that the extra distortion introduced by the coupling step asymptotically vanishes.

Lemma 5. We have

$$\lim_{n \rightarrow \infty} \hat{T}_n(\psi_n^{(n)}, \psi^n) = 0$$

and consequently

$$\limsup_{n \rightarrow \infty} E[\rho_n(X^n, Y^n)] \leq D.$$

In summary, we have shown that there exists a sequence of Model 2 randomized quantizers having rate $R = I_m(\mu \| \psi, D) + \varepsilon$ and asymptotic distortion upper bounded by D which satisfy the output distribution constraint $Y^n \sim \psi^n$. Since $\varepsilon > 0$ is arbitrary, this completes the proof of the achievability of $I_m(\mu \| \psi, D)$ (and the achievability of $D(\mu, \psi, R)$) for finite source and reproduction alphabets.

Remark 6. We note that an obvious approach to achievability would be to generate a codebook where the codewords have i.i.d. components drawn according to ψ . However, the output distribution of the resulting the scheme would be *too far* from the desired ψ^n . In particular, such a scheme produces output \hat{X}^n whose empirical distribution (type) converges to a “favorite type” which is typically different from ψ [23, Theorem 4]. As well, the rate achievable with this scheme at distortion level D is [26, Theorem 2]

$$R = \min_{\psi' \in \mathcal{P}(\mathsf{Y})} (I_m(\mu \| \psi', D) + \mathcal{D}(\psi' \| \psi))$$

which is typically strictly less than $I_m(\mu \| \psi, D)$.

Now let $\mathsf{X} = \mathsf{Y} = \mathbb{R}$, $\rho(x, y) = (x - y)^2$, and assume that μ and ψ have finite second moments. We make use of the final alphabet case to prove achievability for this continuous case. The following lemma provides the necessary link between the two cases.

Lemma 6. There exist a sequence $\{A_k\}$ of finite subsets of \mathbb{R} and sequences of probability measures $\{\mu_k\}$ and $\{\psi_k\}$, both supported on A_k , such that

- (i) $\hat{T}_1(\mu, \mu_k) \rightarrow 0$, $\hat{T}_1(\psi, \psi_k) \rightarrow 0$ as $k \rightarrow \infty$;
- (ii) For any $\varepsilon > 0$ and $D \geq 0$ such that $I_m(\mu \| \psi, D) < \infty$, we have $I_m(\mu_k \| \psi_k, D + \varepsilon) \leq I_m(\mu \| \psi, D)$ for all k large enough.

Let μ_k^n and ψ_k^n denote the n -fold products of μ_k and ψ_k , respectively. Definition (9) of optimal coupling implies that

$\hat{T}_n(\mu^n, \mu_k^n) \leq \hat{T}_1(\mu, \mu_k)$ and $\hat{T}_n(\psi^n, \psi_k^n) \leq \hat{T}_1(\psi, \psi_k)$. Hence for any given $\varepsilon > 0$ by Lemma 6 we can choose k large enough such that for all n ,

$$\hat{T}_n(\mu^n, \mu_k^n) \leq \varepsilon \text{ and } \hat{T}_n(\psi^n, \psi_k^n) \leq \varepsilon, \quad (10)$$

and also $I_m(\mu_k \|\psi_k, D + \varepsilon) \leq I_m(\mu \|\psi, D)$.

Now for each n define the following randomized quantizer:

- (a) Realize the optimal coupling between μ^n and μ_k^n .
- (b) Apply the randomized quantizer scheme for the finite alphabet case with common source and output alphabet A_k , source distribution μ_k^n , and output distribution ψ_k^n . Set the rate of the quantizer to $R = I_m(\mu \|\psi, D) + \varepsilon$.
- (c) Realize the optimal coupling between ψ_k^n and ψ^n .

In particular, the optimal couplings are realized as follows: in (a) the source $X^n \sim \mu^n$ is mapped to $X^n(k) \sim \mu_k^n$, which serves as the source in (b), via $X^n(k) = \hat{f}_{n,k}(X^n, Z_3)$, and in (c) the output $Y^n(k) \sim \psi_k^n$ of the scheme in (b) is mapped to $Y^n \sim \psi^n$ via $Y^n = \hat{g}_{n,k}(Y^n(k), Z_4)$, where Z_3 and Z_4 are uniform randomization variables that are independent of X^n . Thus the composition of these three steps is a valid Model 2 randomized quantizer.

Since $R = I_m(\mu \|\psi, D) + \varepsilon$, in step (b) the asymptotic (in n) distortion $D + \varepsilon$ can be achieved by Lemma 6(ii). Using (10) and the triangle inequality for the norm $\|V^n\|_2 := (\sum_{i=1}^n E[V_i^2])^{1/2}$ on \mathbb{R}^n -valued random vectors having finite second moments, it is straightforward to show that the asymptotic distortion of the overall scheme is upper bounded by $D + l(\varepsilon)$, where $l(\varepsilon) \rightarrow 0$ as $\varepsilon \rightarrow 0$. Since $\varepsilon > 0$ can be taken to be arbitrarily small by choosing k large enough, this completes the achievability proof for the case $X = Y = \mathbb{R}$. \square

VI. CONCLUSION

We investigated a general abstract model for randomized quantization that provides a more suitable framework for certain optimal quantization problems than the ones usually considered in the source coding literature. In particular, our model formalizes the notion of randomly picking a quantizer from the set of all *all* quantizers with a given number of output levels. Using this model, we proved the existence of an optimal randomized vector quantizer under the constraint that the quantizer output has a given distribution.

Our results are mostly non-constructive and it is an open problem how to find (or well approximate) such optimal quantizers. A special case where a scalar source has a density and the output distribution is constrained to be equal to the source distribution was considered in [10] and construction based on dithered uniform quantization followed by a nonlinear mapping was given. Although this construction is optimal in the limit of high resolution ($M \rightarrow \infty$), it is very likely suboptimal for any finite M . In general, it would be interesting to better characterize optimal randomized quantizers in Theorem 3, for example, by finding useful necessary conditions for optimality. It would also be interesting to characterize the high-resolution behavior of the distortion, which should be markedly different from the classical case if the input and output distributions are not equal. Connections between the output distribution-constrained lossy source codes studied in Section V and the

empirical distribution of good rate-distortion codes (see, e.g., [27] and references therein) are also worth studying. Finally, a rigorous theory of randomized quantization paves the way for interesting applications in signaling games in game theory [28] and in stochastic networked control (see [29] and [16] for applications of randomized quantization in real-time coding, and [17] and [30] for quantizers and stochastic kernels viewed as information structures in networked control).

APPENDIX

A. Proof of Lemma 1

For a fixed probability measure μ on X define

$$\Delta_\mu = \{v \in \mathcal{P}(X \times Y) : v(\cdot \times Y) = \mu\}$$

(Δ_μ is the set of all probability measures in $\mathcal{P}(X \times Y)$ whose X -marginal is μ). The following proposition, due to Borkar [13, Lemma 2.2], gives a characterization of the extreme points of Δ_μ .

Proposition 1. Δ_μ is closed and convex, and its set of extreme points $\Delta_{\mu,e}$ is a Borel set in $\mathcal{P}(X \times Y)$. Furthermore, $v \in \Delta_{\mu,e}$ if and only if $v(dx dy)$ can be disintegrated as

$$v(dx dy) = Q(dy|x)\mu(dx)$$

where $Q(\cdot|x)$ is a Dirac measure for μ -a.e. x , i.e., there exists a measurable function $f : X \rightarrow Y$ such that $Q(\cdot|x) = \delta_{f(x)}(\cdot)$ for μ -a.e. x .

In fact, Borkar did not explicitly state Borel measurability of $\Delta_{\mu,e}$ in [13], but the proof of [13, Lemma 2.3] clearly implies this.

By Proposition 1 it is clear that $v \in \Gamma_\mu(M)$ if and only if $v \in \Delta_{\mu,e}$ and its marginal on Y is supported on a set having at most M elements, i.e., for some $L \leq M$ and $\{y_1, \dots, y_L\} \subset Y$,

$$v(X \times \{y_1, \dots, y_L\}) = 1.$$

Let $\{y_n\}_{n \geq 1}$ be a countable dense subset of Y and define following subsets of $\Delta_{\mu,e}$:

$$\Omega_k = \bigcup_{n_1 \geq 1, \dots, n_M \geq 1} \left\{ v \in \Delta_{\mu,e} : v\left(X \times \bigcup_{i=1}^M B(y_{n_i}, 1/k)\right) = 1 \right\}$$

and

$$\Sigma = \bigcap_{k=1}^{\infty} \Omega_k$$

where $B(y, r)$ denotes the open ball in Y centered at y having radius r . Sets of the form

$$\left\{ v \in \mathcal{P}(X \times Y) : v\left(X \times \bigcup_{i=1}^M B(y_{n_i}, 1/k)\right) = 1 \right\}$$

are Borel sets by [31, Proposition 7.25]. Since $\Delta_{\mu,e}$ is a Borel set, Ω_k is a Borel set for all k . Thus Σ is a Borel set in $\mathcal{P}(X \times Y)$. We will prove that $\Sigma = \Gamma_\mu(M)$.

Since $\{y_n\}_{n \geq 1}$ is dense in Y , for any $v \in \Gamma_\mu(M)$ and $k \geq 1$ there exist $\tilde{n}_1, \dots, \tilde{n}_M$ such that $\text{supp}(v(X \times \cdot)) \subset$

$\bigcup_{i=1}^M B(y_{\tilde{n}_i}, 1/k)$. Thus $\Gamma_\mu(M) \subset \Omega_k$ for all k , implying $\Gamma_\mu(M) \subset \Sigma$.

To prove the inclusion $\Sigma \subset \Gamma_\mu(M)$, let $v \in \Sigma$ and notice that for all k there exist $n_1^k, n_2^k, \dots, n_M^k$ such that

$$v\left(X \times \bigcup_{i=1}^M B(y_{n_i^k}, 1/k)\right) = 1.$$

Let us define $K_n = X \times \bigcap_{k=1}^n \bigcup_{i=1}^M B(y_{n_i^k}, 1/k)$. Clearly, $K_{n+1} \subset K_n$ and $v(K_n) = 1$, for all n . Letting

$$G = \bigcap_{k=1}^{\infty} \bigcup_{i=1}^M B(y_{n_i^k}, 1/k),$$

we have $v(X \times G) = 1$. If we can prove that G has at most M distinct elements, then $v \in \Gamma_\mu(M)$. Assuming the contrary, there must exist distinct $\{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_M, \hat{y}_{M+1}\} \subset G$. Let $\varepsilon = \min\{\|\hat{y}_i - \hat{y}_j\| : i, j = 1, \dots, M+1, i \neq j\}$. Clearly, for $\frac{1}{k} < \frac{\varepsilon}{4}$, $\bigcup_{i=1}^M B(y_{n_i^k}, 1/k)$ cannot contain $\{\hat{y}_1, \hat{y}_2, \dots, \hat{y}_M, \hat{y}_{M+1}\}$, a contradiction. Thus G has at most M elements and we obtain $\Sigma = \Gamma_\mu(M)$. \square

B. Proof of Theorem 1

We will need the following result which gives a necessary and sufficient condition for the measurability of a mapping from a measurable space to $\mathcal{P}(E)$, where E is a Polish space. It is proved for compact E in [32, Theorem 2.1] and for noncompact E it is the corollary of [31, Proposition 7.25].

Theorem 8. *Let (Ω, \mathcal{F}) be a measurable space and let E be a Polish space. A mapping $h : \Omega \rightarrow \mathcal{P}(E)$ is measurable if and only if the real valued functions $\omega \mapsto h(\omega)(A)$ from Ω to $[0, 1]$ are measurable for all $A \in \mathcal{B}(E)$.*

For any (q, ν) define $f : \mathbb{R}^m \rightarrow \Gamma_\mu(M)$ by $f(z) = \delta_{q(x,z)}(dy)\mu(dx)$. By Theorem 8, f is measurable if and only if the mappings $z \mapsto \int \delta_{q(x,z)}(C_x)\mu(dx)$ are measurable for all $C \in \mathcal{B}(X \times Y)$, where $C_x = \{y : (x, y) \in C\}$. Observe that $\delta_{q(x,z)}(C_x)$ is a measurable function of (x, z) because $\{(x, z) \in X \times Z : \delta_{q(x,z)}(C_x) = 1\} = \{(x, z) \in X \times Z : (x, q(x, z)) \in C\}$. By [33, Theorem 18.3] $\int \delta_{q(x,z)}(C_x)\mu(dx)$ is measurable as well. Hence f is measurable.

Thus we can define the probability measure P supported on $\Gamma_\mu(M)$ by $P = \nu \circ f^{-1}$ (i.e., $P(B) = \nu(f^{-1}(B))$ for any Borel set $B \subset \Gamma_\mu(M)$). Then, for the corresponding v_P we have $(X, Y) \sim v_P$, i.e., for $C \in \mathcal{B}(X \times Y)$,

$$\begin{aligned} \Pr\{(X, q(X, Z)) \in C\} &= \int_Z \int_X \delta_{q(x,z)}(C_x)\mu(dx)\nu(dz) \\ &= \int_Z f(z)(C)\nu(dz) \\ &= \int_{\Gamma_\mu(M)} v(C)P(dv) \\ &= v_P(C). \end{aligned}$$

Conversely, let v_P be defined as in (3) with P supported on $\Gamma_\mu(M)$, i.e., $v_P = \int_{\Gamma_\mu(M)} vP(dv)$. Define the mapping $\Gamma_\mu(M) \ni v \mapsto q_v$, where q_v is the μ -a.e. defined quantizer in \mathcal{Q}_M , giving $v(dx dy) = \mu(dx)\delta_{q_v(x)}(dy)$. Since $\Gamma_\mu(M)$ is

an uncountable Borel space, there is a measurable bijection (Borel isomorphism) $g : \mathbb{R}^m \rightarrow \Gamma_\mu(M)$ between \mathbb{R}^m and $\Gamma_\mu(M)$ [21]. Now define q by $q(x, z) = q_{g(z)}(x)$ and let $\nu = P \circ g$. Then for all z , $q(\cdot, z)$ is a μ -a.e. defined M -level quantizer. However, it is not clear whether $q(x, z)$ is measurable. Therefore we will construct another measurable function $\tilde{q}(x, z)$ such that $\tilde{q}(\cdot, z)$ is an M -level quantizer and $\tilde{q}(\cdot, z) = q(\cdot, z)$ μ -a.e., for all z . Then we will prove that $(X, Y) = (X, \tilde{q}(X, Z)) \sim v_P$ where $Z \sim \nu$. Define the stochastic kernel on $X \times Y$ given $\Gamma_\mu(M)$ as

$$\gamma(dx dy|v) = v(dx dy).$$

Clearly, γ is well defined because $\Gamma_\mu(M)$ is a Borel subset of $\mathcal{P}(X \times Y)$. Observe that for each $v \in \Gamma_\mu(M)$, we have

$$\gamma(C|v) = \int_X \delta_{q_v(x)}(C_x)\mu(dx) \quad (11)$$

for $C \in \mathcal{B}(X \times Y)$. Furthermore, by [31, Proposition 7.27] there exists a stochastic kernel $\eta(dy|x, v)$ on Y given $X \times \Gamma_\mu(M)$ which satisfies for all $C \in \mathcal{B}(X \times Y)$ and $v \in \Gamma_\mu(M)$,

$$\gamma(C|v) = \int_X \eta(C_x|x, v)\mu(dx). \quad (12)$$

Since $\mathcal{B}(Y)$ is countably generated by the separability of Y , for any $v \in \Gamma_\mu(M)$ we have $\eta(\cdot|x, v) = \delta_{q_v(x)}(\cdot)$ μ -a.e. by (11) and (12). Since η is a stochastic kernel, it can be represented as a measurable function from $X \times \Gamma_\mu(M)$ to $\mathcal{P}(Y)$, i.e.,

$$\eta : X \times \Gamma_\mu(M) \rightarrow \mathcal{P}(Y).$$

Define $\mathcal{P}_1(Y) = \{\psi \in \mathcal{P}(Y) : \psi(\{y\}) = 1 \text{ for some } y \in Y\}$. $\mathcal{P}_1(Y)$ is a closed (thus measurable) subset of $\mathcal{P}(Y)$ by [34, Lemma 6.2]. Hence, $M := \eta^{-1}(\mathcal{P}_1(Y))$ is also measurable. Observe that for any $v \in \Gamma_\mu(M)$ we have $M_v := \{x \in X : (x, v) \in M\} \supset \{x \in X : \eta(\cdot|x, v) = \delta_{q_v(x)}(\cdot)\}$. Thus $\mu(M_v) = 1$ for all $v \in \Gamma_\mu(M)$, which implies $\mu \otimes P(M) = 1$. Define the function \tilde{q}_v from $X \times \Gamma_\mu(M)$ to Y as

$$\tilde{q}_v(x) = \begin{cases} \tilde{y}, & \text{if } (x, v) \in M, \text{ where } \eta(\{\tilde{y}\}|x, v) = 1, \\ y, & \text{otherwise,} \end{cases}$$

where y is fixed. By construction, $\tilde{q}_v(x) = q_v(x)$ μ -a.e., for all $v \in \Gamma_\mu(M)$. For any $C \in \mathcal{B}(Y)$ we have

$$\begin{aligned} \tilde{q}_v^{-1}(C) &= \{(x, v) \in X \times \Gamma_\mu(M) : \tilde{q}_v(x) \in C\} \\ &= \{(x, v) \in M : \tilde{q}_v(x) \in C\} \cup \{(x, v) \in M^c : \tilde{q}_v(x) \in C\}. \end{aligned}$$

Clearly $\{(x, v) \in M^c : \tilde{q}_v(x) \in C\} = M^c$ or \emptyset depending on whether or not y is an element of C . Hence, $\tilde{q}_v^{-1}(C) \in \mathcal{B}(X \times \Gamma_\mu(M))$ if $\{(x, v) \in M : \tilde{q}_v(x) \in C\} \in \mathcal{B}(X \times \Gamma_\mu(M))$. But $\{(x, v) \in M : \tilde{q}_v(x) \in C\} = \{(x, v) \in M : \eta(C|x, v) = 1\}$ which is in $\mathcal{B}(X \times \Gamma_\mu(M))$ by the measurability of $\eta(C|\cdot, \cdot)$. Thus, \tilde{q} is a measurable function from $X \times \Gamma_\mu(M)$ to Y .

Let us define \tilde{q} as $\tilde{q}(x, z) = \tilde{q}_{g(z)}(x)$. By the measurability of g it is clear that \tilde{q} is measurable. In addition, for any $z \in Z$ $\tilde{q}(\cdot, z)$ is an M -level quantizer which is μ -a.e. equal to $q(\cdot, z)$. Finally, if $Z \sim \nu$ is independent of X and $Y = \tilde{q}(X, Z)$, then $(X, Y) \sim v_P$, i.e.,

$$\begin{aligned}
& \Pr\{(X, \tilde{q}(X, Z)) \in C\} \\
&= \int_{\mathbb{Z}} \int_{\mathbb{X}} \delta_{\tilde{q}(x,z)}(C_x) \mu(dx) \nu(dz) \\
&= \int_{\Gamma_\mu(M)} \int_{\mathbb{X}} \delta_{\tilde{q}_v(x)}(C_x) \mu(dx) P(dv) \\
&= \int_{\Gamma_\mu(M)} \int_{\mathbb{M}_v} \eta(C_x|x, v) \mu(dx) P(dv) \\
&= \int_{\Gamma_\mu(M)} \gamma(C|v) P(dv) \\
&= \int_{\Gamma_\mu(M)} v(C) P(dv) \\
&= v_p(C). \quad \square
\end{aligned}$$

C. Proof of Theorem 2

If (e, d, ν) is a Model 1 randomized quantizer, then setting $q(x, z) = d(e(x, z), z)$ defines a Model 2 randomized quantizer (q, ν) such that the joint distributions of their inputs and outputs coincide.

Conversely, let (q, ν) be a Model 2 randomized quantizer. It is obvious that q can be decomposed into an encoder $e : \mathbb{X} \times \mathbb{Z} \rightarrow \{1, \dots, M\}$ and decoder $d : \{1, \dots, M\} \times \mathbb{Z} \rightarrow \mathbb{Y}$ such that $d(e(x, z), z) = q(x, z)$ for all x and z . The difficulty lies in showing that this can be done so that the resulting e and d are measurable. In fact, we instead construct measurable e and d whose composition is $\mu \otimes \nu$ -a.e. equal to q , which is sufficient to imply the theorem.

Let (q, ν) be a Model 2 randomized quantizer. Since \mathbb{R}^n and $[0, 1]$ are both uncountable Borel spaces, there exists a Borel isomorphism $f : \mathbb{R}^n \rightarrow [0, 1]$ [21]. Define $\hat{q} : \mathbb{X} \times \mathbb{Y} \rightarrow [0, 1]$ by $\hat{q} = f \circ q$. Hence, \hat{q} is measurable and, for any fixed z , $\hat{q}(\cdot, z)$ is an M -level quantizer from \mathbb{X} to $[0, 1]$. Also note that $q = f^{-1} \circ \hat{q}$.

Now for any fixed $z \in \mathbb{Z}$ consider only those output points of $\hat{q}(\cdot, z)$ that occur with *positive* μ probability and order these according to their magnitude from the smallest to the largest. For $i = 1, \dots, M$ let the function $f_i(z)$ take the value of the i th smallest such output point. If there is no such value, let $f_i(z) = 1$. We first prove that all the f_i are measurable and then define the encoder and the decoder in terms of these functions.

Observe that for any $a \in [0, 1]$, by definition

$$\{z \in \mathbb{Z} : f_1(z) \leq a\} = \left\{z \in \mathbb{Z} : \int_{\mathbb{X}} \delta_{\hat{q}(x,z)}([0, a]) \mu(dx) > 0\right\},$$

where the set on the right hand side is a Borel set by Fubini's theorem. Hence, f_1 is a measurable function. Define $E_1 = \{(x, z) \in \mathbb{X} \times \mathbb{Z} : \hat{q}(x, z) - f_1(z) = 0\}$, a Borel set. Letting $E_{1,z} = \{x \in \mathbb{X} : (x, z) \in E_1\}$ denote the z -section of E_1 , for any $a \in [0, 1]$ we have

$$\begin{aligned}
& \{z \in \mathbb{Z} : f_2(z) \leq a\} \\
&= \left\{z \in \mathbb{Z} : \int_{\mathbb{X} \setminus E_{1,z}} \delta_{\hat{q}(x,z)}([0, a]) \mu(dx) > 0\right\},
\end{aligned}$$

and thus f_2 is measurable. Continuing in this fashion, we define the Borel sets $E_i = \{(x, z) : \hat{q}(x, z) - f_i(z) = 0\}$ and write, for any $a \in [0, 1]$,

$$\begin{aligned}
& \{z \in \mathbb{Z} : f_i(z) \leq a\} \\
&= \left\{z \in \mathbb{Z} : \int_{\mathbb{X} \setminus \bigcup_{j=1}^{i-1} E_{j,z}} \delta_{\hat{q}(x,z)}([0, a]) \mu(dx) > 0\right\},
\end{aligned}$$

proving that f_i is measurable for all $i = 1, \dots, M$.

Define

$$\begin{aligned}
N &= \{(x, z) \in \mathbb{X} \times \mathbb{Z} : \hat{q}(x, z) \neq f_i(z) \text{ for all } i = 1, \dots, M\} \\
&= \mathbb{X} \times \mathbb{Z} \setminus \bigcup_{i=1}^M E_i.
\end{aligned}$$

Clearly, N is a Borel set and $\mu \otimes \nu(N) = 0$ by Fubini's theorem and the definition of f_1, \dots, f_M . Now we can define

$$e(x, z) = \sum_{i=1}^M i 1_{\{\hat{q}(x,z)=f_i(z)\}} + M 1_N(x, z)$$

and

$$d(i, z) = \sum_{j=1}^M f^{-1} \circ f_j(z) 1_{\{i=j\}},$$

where 1_B denotes the indicator of event (or set) B . The measurability of \hat{q} and f, f_1, \dots, f_M implies that e and d are measurable. Since $d(e(x, z), z) = \hat{q}(x, z)$ $\mu \otimes \nu$ -a.e. by construction, this completes the proof. \square

D. Proof of Theorem 3

1) Proof under Assumption 1

To simplify the notation we redefine the reconstruction alphabet as $\mathbb{Y} = B$, so that \mathbb{Y} is a compact subset of \mathbb{R}^n . It follows from the continuity of ρ that L is lower semicontinuous on $\mathcal{P}(\mathbb{X} \times \mathbb{Y})$ for the weak topology (see, e.g., [22, Lemma 4.3]). Hence, to show the existence of a minimizer for problem **(P1)** it would suffice to prove that $\Gamma_{\mu, \psi}^{\mathbb{R}}(M) = \Gamma_{\mu}^{\mathbb{R}}(M) \cap \Gamma_{\mu, \psi}$ is compact. It is known that $\Gamma_{\mu, \psi}$ is compact [22, Chapter 4], but unfortunately $\Gamma_{\mu}^{\mathbb{R}}(M)$ is not closed [17] and it seems doubtful that $\Gamma_{\mu}^{\mathbb{R}}(M)$ is compact. Hence, we will develop a different argument which is based on optimal transport theory. We will first give the proof under Assumption 1; the proof under Assumption 2 then follows via a one-point compactification argument.

Let $\mathcal{P}_M(\mathbb{Y}) = \{\psi_0 \in \mathcal{P}(\mathbb{Y}) : |\text{supp}(\psi_0)| \leq M\}$ be the set of discrete distributions with M atoms or less on \mathbb{Y} .

Lemma 7. $\mathcal{P}_M(\mathbb{Y})$ is compact in $\mathcal{P}(\mathbb{Y})$.

Proof: Let $\{\psi_n\}$ be an arbitrary sequence in $\mathcal{P}_M(\mathbb{Y})$. Each ψ_n can be represented by points $(y_1^n, \dots, y_M^n) = y^n \in \mathbb{Y}^M$ and $(p_1^n, \dots, p_M^n) = p^n \in K_s$, where $K_s = \{(p_1, \dots, p_M) \in \mathbb{R}^M : \sum_{i=1}^M p_i = 1, p_i \geq 0\}$ is the probability simplex in \mathbb{R}^M . Let $w_n = (y^n, p^n)$. Since $\mathbb{Y}^M \times K_s$ is compact, there exists a subsequence $\{w^{n_k}\}$ converging to some w in $\mathbb{Y}^M \times K_s$. Let ψ be the probability measure in $\mathcal{P}_M(\mathbb{Y})$ which is represented by w . It straightforward to show that ψ is a weak limit of $\{\psi^{n_k}\}$. This completes the proof. \square

Define

$$\hat{\Gamma}_\mu(M) = \bigcup_{\psi_0 \in \mathcal{P}_M(Y)} \{ \hat{v} \in \Gamma_{\mu, \psi_0} : L(\hat{v}) = \min_{v \in \Gamma_{\mu, \psi_0}} L(v) \}.$$

The elements of $\hat{\Gamma}_\mu(M)$ are the probability measures which solve the optimal transport problem (see, e.g., [22]) for fixed input marginal μ and some output marginal ψ_0 in $\mathcal{P}_M(Y)$. At the end of this proof Lemma 11 shows that $\hat{\Gamma}_\mu(M)$ is a Borel set. Let $\hat{\Gamma}_\mu^R(M)$ be the randomization of $\hat{\Gamma}_\mu(M)$, obtained by replacing $\Gamma_\mu(M)$ with $\hat{\Gamma}_\mu(M)$ in (4). Define the optimization problem **(P2)** as

$$\begin{aligned} \text{(P2) minimize } & L(v) \\ \text{subject to } & v \in \hat{\Gamma}_{\mu, \psi}^R(M), \end{aligned}$$

where $\hat{\Gamma}_{\mu, \psi}^R(M) = \hat{\Gamma}_\mu^R(M) \cap \Gamma_{\mu, \psi}$.

Proposition 2. *For any $v^* \in \Gamma_{\mu, \psi}^R(M)$ there exists $\hat{v} \in \hat{\Gamma}_{\mu, \psi}^R(M)$ such that $L(v^*) \geq L(\hat{v})$. Hence, the distortion of any minimizer in **(P2)** is less than or equal to the distortion of a minimizer in **(P1)**.*

To prove Proposition 2 we need the following lemma.

Lemma 8. *Let P be a probability measure on $\Gamma_\mu(M)$. Then there exists a measurable mapping $f : \Gamma_\mu(M) \rightarrow \hat{\Gamma}_\mu(M)$ such that $v(X \times \cdot) = f(v)(X \times \cdot)$ and $L(v) \geq L(f(v))$, P -a.e.*

Proof: Define the projections $f_1 : \Gamma_\mu(M) \rightarrow \mathcal{P}_M(Y)$ and $f_2 : \hat{\Gamma}_\mu(M) \rightarrow \mathcal{P}_M(Y)$ by $f_1(v) = v(X \times \cdot)$, $f_2(v) = v(X \times \cdot)$. Note that f_1 is continuous and f_2 is continuous and onto. Define $\tilde{P} = P \circ f_1^{-1}$ on $\mathcal{P}_M(Y)$. By Yankov's lemma [35, Appendix 3] there exists a mapping g from $\mathcal{P}_M(Y)$ to $\hat{\Gamma}_\mu(M)$ such that $f_2(g(\psi)) = \psi$ \tilde{P} -a.e. Then, it is straightforward to show that $f = g \circ f_1$ satisfies conditions $v(X \times \cdot) = f(v)(X \times \cdot)$ and $L(v) \geq L(f(v))$, P -a.e. \square

Proof of Proposition 2: Let $v^* \in \Gamma_{\mu, \psi}^R(M)$, i.e.,

$$v^* = \int_{\Gamma_\mu(M)} v P(dv) \text{ and } v^*(X \times \cdot) = \psi.$$

By Lemma 8 there exists $f : \Gamma_\mu(M) \rightarrow \hat{\Gamma}_\mu(M)$ such that $v(X \times \cdot) = f(v)(X \times \cdot)$ and $L(v) \geq L(f(v))$, P -a.e. Define $\tilde{P} = P \circ f^{-1} \in \mathcal{P}(\hat{\Gamma}_\mu(M))$ and $\hat{v} = \int_{\hat{\Gamma}_\mu(M)} v \tilde{P}(dv) \in \hat{\Gamma}_\mu^R(M)$. We have

$$\begin{aligned} L(v^*) &= \int_{\Gamma_\mu(M)} L(v) P(dv) \geq \int_{\Gamma_\mu(M)} L(f(v)) P(dv) \\ &= \int_{\hat{\Gamma}_\mu(M)} L(v) \tilde{P}(dv) = L(\hat{v}) \end{aligned}$$

as well as

$$\begin{aligned} v^*(X \times \cdot) &= \int_{\Gamma_\mu(M)} v(X \times \cdot) P(dv) \\ &= \int_{\Gamma_\mu(M)} f(v)(X \times \cdot) P(dv) \\ &= \int_{\hat{\Gamma}_\mu(M)} v(X \times \cdot) \tilde{P}(dv) = \hat{v}(X \times \cdot). \end{aligned}$$

This completes the proof. \square

Recall the set Δ_μ and its set of its extreme points $\Delta_{\mu, e}$ from Proposition 1. It is proved in [13] and [36] that any $\tilde{v} \in \Delta_\mu$ can be written as $\tilde{v} = \int_{\Delta_{\mu, e}} v P(dv)$ for some $P \in \mathcal{P}(\Delta_{\mu, e})$. By Proposition 1 we also have $\Gamma_\mu(M) \subset \Delta_{\mu, e}$. The following lemma is based on these two facts.

Lemma 9. *Let $\tilde{v} \in \Delta_\mu$ which is represented as $\tilde{v} = \int_{\Delta_{\mu, e}} v P(dv)$. If $\tilde{v}(X \times \cdot) \in \mathcal{P}_M(Y)$, then $P(\Gamma_\mu(M)) = 1$.*

Proof: Since $\tilde{v}(X \times \cdot) \in \mathcal{P}_M(Y)$, there exist a finite set $B \subset Y$ having $M' \leq M$ elements such that $\tilde{v}(X \times B) = 1$. We have

$$\begin{aligned} \tilde{v}(X \times B) &= \int_{\Delta_{\mu, e}} v(X \times B) P(dv) \\ &= \int_{\Delta_{\mu, e} \setminus \Gamma_\mu(M)} v(X \times B) P(dv) \\ &\quad + \int_{\Gamma_\mu(M)} v(X \times B) P(dv). \end{aligned}$$

Since $v(X \times B) < 1$ for all $v \in \Delta_{\mu, e} \setminus \Gamma_\mu(M)$, we obtain $P(\Gamma_\mu(M)) = 1$. \square

Lemma 9 implies $\hat{\Gamma}_\mu(M) \subset \Gamma_\mu^R(M)$ because $v(X \times \cdot) \in \mathcal{P}_M(Y)$ when $v \in \hat{\Gamma}_\mu(M)$. Define $h : \mathcal{P}(\Gamma_\mu(M)) \rightarrow \Delta_\mu$ as follows:

$$h(P)(\cdot) = \int_{\Gamma_\mu(M)} v(\cdot) P(dv). \quad (13)$$

It is clear that the range of h is $\Gamma_\mu^R(M) \subset \Delta_\mu$.

Lemma 10. *h is continuous.*

Proof: Assume $\{P_n\}$ converges weakly to P in $\mathcal{P}(\Gamma_\mu(M))$. Then, for any continuous and bounded real function f on $X \times Y$

$$\begin{aligned} \lim_{n \rightarrow \infty} \int_{\Gamma_\mu(M)} \int_{X \times Y} f(x, y) v(dx dy) P_n(dv) \\ = \int_{\Gamma_\mu(M)} \int_{X \times Y} f(x, y) v(dx dy) P(dv) \end{aligned}$$

if the mapping $v \mapsto \int_{X \times Y} f(x, y) v(dx dy)$ is continuous and bounded on $\Gamma_\mu(M)$. Clearly this mapping is continuous by the definition of weak convergence and bounded by the boundedness of f . Thus

$$\int_{\Gamma_\mu(M)} v P_n(dv) \rightarrow \int_{\Gamma_\mu(M)} v P(dv)$$

weakly, completing the proof. \square

Since $\hat{\Gamma}_\mu(M) \subset \Gamma_\mu^R(M)$, we have $\mathcal{P}^{\text{opt}}(\Gamma_\mu(M)) := h^{-1}(\hat{\Gamma}_\mu(M)) \subset \mathcal{P}(\Gamma_\mu(M))$, which is measurable by the measurability of $\hat{\Gamma}_\mu(M)$ and h . Let $g : \mathcal{P}^{\text{opt}}(\Gamma_\mu(M)) \rightarrow \hat{\Gamma}_\mu(M)$ be the restriction of h to $\mathcal{P}^{\text{opt}}(\Gamma_\mu(M))$. Clearly g is measurable and onto. By Yankov's lemma [35] for any probability measure P on $\hat{\Gamma}_\mu(M)$ there exists a measurable mapping $\varphi : \hat{\Gamma}_\mu(M) \rightarrow \mathcal{P}^{\text{opt}}(\Gamma_\mu(M))$ such that $g(\varphi(\hat{v})) = \hat{v}$ P -a.e. In addition, since $\varphi(\hat{v}) \in g^{-1}(\hat{v})$ P -a.e., we have

$$L(\hat{v}) = \int_{\Gamma_\mu(M)} L(v) \varphi(\hat{v})(dv) \quad (14)$$

and

$$\hat{v}(\mathbb{X} \times \cdot) = \int_{\Gamma_\mu(M)} v(\mathbb{X} \times \cdot) \varphi(\hat{v})(dv) \quad (15)$$

P -a.e. Define the stochastic kernel $\Pi(dv|\hat{v})$ on $\Gamma_\mu(M)$ given $\hat{\Gamma}_\mu(M)$ as

$$\Pi(dv|\hat{v}) = \varphi(\hat{v})(dv). \quad (16)$$

Since φ is measurable, $\Pi(dv|\hat{v})$ is well defined. Observe that both φ and $\Pi(dv|\hat{v})$ depend on the probability measure $P \in \hat{\Gamma}_\mu(M)$.

Proposition 3. *If (P2) has a minimizer v^* , then we can find $\bar{v} \in \Gamma_{\mu,\psi}^{\text{R}}(M)$ such that $L(\bar{v}) = L(v^*)$, implying that \bar{v} is a minimizer for (P1).*

Proof: v^* can be written as $v^* = \int_{\hat{\Gamma}_\mu(M)} \hat{v} P(d\hat{v})$. Consider the stochastic kernel $\Pi(dv|\hat{v})$ defined in (16). Composing P with Π we obtain a probability measure Λ on $\hat{\Gamma}_\mu(M) \times \Gamma_\mu(M)$ given by

$$\Lambda(d\hat{v} dv) = P(d\hat{v})\Pi(dv|\hat{v}). \quad (17)$$

Let $\tilde{P} = \Lambda(\hat{\Gamma}_\mu(M) \times \cdot) \in \mathcal{P}(\Gamma_\mu(M))$. Define the randomized quantizer $\bar{v} \in \Gamma_\mu^{\text{R}}(M)$ as $\bar{v} = \int_{\Gamma_\mu(M)} v \tilde{P}(dv)$. We show that $L(v^*) = L(\bar{v})$ and $v^*(\mathbb{X} \times \cdot) = \bar{v}(\mathbb{X} \times \cdot)$ which will complete the proof. We have

$$\begin{aligned} L(v^*) &= \int_{\hat{\Gamma}_\mu(M)} L(\hat{v}) P(d\hat{v}) \\ &= \int_{\hat{\Gamma}_\mu(M)} \int_{\Gamma_\mu(M)} L(v) \varphi(\hat{v})(dv) P(d\hat{v}) \quad (\text{by (14)}) \\ &= \int_{\hat{\Gamma}_\mu(M) \times \Gamma_\mu(M)} L(v) \Lambda(d\hat{v} dv) \quad (\text{by (16)}) \\ &= \int_{\Gamma_\mu(M)} L(v) \tilde{P}(dv) = L(\bar{v}). \end{aligned}$$

Similarly,

$$\begin{aligned} v^*(\mathbb{X} \times \cdot) &= \int_{\hat{\Gamma}_\mu(M)} \hat{v}(\mathbb{X} \times \cdot) P(d\hat{v}) \\ &= \int_{\hat{\Gamma}_\mu(M)} \int_{\Gamma_\mu(M)} v(\mathbb{X} \times \cdot) \varphi(\hat{v})(dv) P(d\hat{v}) \quad (\text{by (15)}) \\ &= \int_{\hat{\Gamma}_\mu(M) \times \Gamma_\mu(M)} v(\mathbb{X} \times \cdot) \Lambda(d\hat{v} dv) \quad (\text{by (16)}) \\ &= \int_{\Gamma_\mu(M)} v(\mathbb{X} \times \cdot) \tilde{P}(dv) = \bar{v}(\mathbb{X} \times \cdot). \end{aligned}$$

By Proposition 2, \bar{v} is a minimizer for (P1). \square

Hence, to prove the existence of a minimizer for (P1) it is enough to prove the existence of a minimizer for (P2). Before proceeding to the proof we need to define the optimal transport problem. Optimal transport problem for marginals $\pi \in \mathcal{P}(X)$, $\lambda \in \mathcal{P}(Y)$ and cost function $c : X \times Y \rightarrow [0, \infty]$ is defined as:

$$\begin{aligned} &\text{minimize } \int_{X \times Y} c(x, y) v(dx dy) \\ &\text{subject to } v \in \Gamma_{\pi, \lambda}. \end{aligned} \quad (18)$$

The following result is about the structure of the optimal v in (18). It uses the concept of c -cyclically monotone sets [22,

Definition 5.1]. A set $B \subset X \times Y$ is said to be c -cyclically monotone if for any $N \geq 1$ and pairs $(x_1, y_1), \dots, (x_N, y_N)$ in B , the following inequality holds:

$$\sum_{i=1}^N c(x_i, y_i) \leq \sum_{i=1}^N c(x_i, y_{i+1}),$$

where $y_{N+1} := y_1$.

Informally, when $v \in \Gamma_{\pi, \lambda}$ is concentrated on a c -cyclically monotone set, then its cost cannot be improved by local perturbations; see the discussion in [22, Chapter 5]. The following result shows that an optimal v must concentrate on a c -cyclically monotone set.

Proposition 4 ([37, Theorem 1.2], [22, Theorem 5.10]). *Let $c : X \times Y \rightarrow [0, \infty]$ be continuous. If $v \in \Gamma_{\pi, \lambda}$ is a solution to the optimal transport problem (18) and $\int_{X \times Y} c(x, y) v(dx dy) < \infty$, then v is concentrated on some c -cyclically monotone set.*

For any $K \subset \mathcal{P}(X)$ and $S \subset \mathcal{P}(Y)$ define $\Xi_{K, S} \subset \mathcal{P}(X \times Y)$ as the set of probability measures which are concentrated on some c -cyclically monotone set and solve (18) for some $\pi \in K$, $\lambda \in S$. The following result is a slight modification of [22, Corollary 5.21].

Proposition 5. *If K and S are compact, then $\Xi_{K, S}$ is compact.*

Proof: Let $\{v_n\}$ be a sequence in $\Xi_{K, S}$. It can be shown that there exists a subsequence $\{v_{n_k}\}$ converging to v whose marginals belong to K and S [22, Lemma 4.4]. Since each v_{n_k} is concentrated on a c -cyclically monotone set by assumption, it can be shown by using the continuity of c that v is also concentrated on a c -cyclically monotone set (see proof of Theorem 5.20 in [22]). Then v is also an element of $\Xi_{K, S}$ by [37, Theorem B]. \square

Since $\{\mu\}$ and $\mathcal{P}_M(Y)$ are both compact, we obtain that $\Xi_{\{\mu\}, \mathcal{P}_M(Y)}$ is compact. Thus it follows that $\mathcal{P}(\Xi_{\{\mu\}, \mathcal{P}_M(Y)})$ is also compact. Furthermore, by Proposition 4 we have $\Xi_{\{\mu\}, \mathcal{P}_M(Y)} \supset \{v \in \hat{\Gamma}_\mu(M) : L(v) < \infty\}$. Hence the randomization can be restricted to $\Xi_{\{\mu\}, \mathcal{P}_M(Y)}$ when defining $\hat{\Gamma}_\mu^{\text{R}}(M)$ for (P2). Let $\Xi_{\{\mu\}, \mathcal{P}_M(Y)}^{\text{R}}$ be the randomization of $\Xi_{\{\mu\}, \mathcal{P}_M(Y)}$ obtained by replacing $\Gamma_\mu(M)$ with $\Xi_{\{\mu\}, \mathcal{P}_M(Y)}$ in (4). One can show that the mapping $\mathcal{P}(\Xi_{\{\mu\}, \mathcal{P}_M(Y)}) \ni P \mapsto v_P \in \Xi_{\{\mu\}, \mathcal{P}_M(Y)}^{\text{R}}$ is continuous by using the same proof as in Lemma 10. Thus $\Xi_{\{\mu\}, \mathcal{P}_M(Y)}^{\text{R}}$ is the continuous image of a compact set, and thus it is also compact. This, together with the compactness of $\Gamma_{\mu, \psi}$ and the lower semicontinuity of L , implies the existence of the minimizer for (P2) under Assumption 1.

To tie up a loose end, we still have to show that $\hat{\Gamma}_\mu(M)$ is measurable, which will complete the proof under Assumption 1.

Lemma 11. *$\hat{\Gamma}_\mu(M)$ is a Borel set.*

Proof: Let us define $\hat{\Gamma}_\mu^f(M) := \{v \in \hat{\Gamma}_\mu(M) : L(v) < \infty\}$ and $\hat{\Gamma}_\mu^\infty(M) = \hat{\Gamma}_\mu(M) \setminus \hat{\Gamma}_\mu^f(M)$. Since solutions to the optimal transport problem having finite costs must concentrate on c -cyclically monotone sets by Proposition 4, we

have $\hat{\Gamma}_\mu^f(M) = \{v \in \Xi_{\{\mu\}, \mathcal{P}_M(\mathcal{Y})} : L(v) < \infty\}$. Hence, $\hat{\Gamma}_\mu^f(M)$ is a Borel set since $\Xi_{\{\mu\}, \mathcal{P}_M(\mathcal{Y})}$ is compact and L is lower semi-continuous. Recall the continuous mapping f_2 in the proof of Lemma 8. Since $\Xi_{\{\mu\}, \mathcal{P}_M(\mathcal{Y})}$ is compact, $\{v \in \Xi_{\{\mu\}, \mathcal{P}_M(\mathcal{Y})} : L(v) \leq N\}$ is also compact for all $N \geq 0$. Hence, $f_2(\hat{\Gamma}_\mu^f(M)) = \bigcup_{N=0}^{\infty} f_2(\{v \in \Xi_{\{\mu\}, \mathcal{P}_M(\mathcal{Y})} : L(v) \leq N\})$ is σ -compact, so a Borel set, in $\mathcal{P}_M(\mathcal{Y})$. Since $f_2(\hat{\Gamma}_\mu^\infty(M)) = \mathcal{P}_M(\mathcal{Y}) \setminus f_2(\hat{\Gamma}_\mu^f(M))$, $f_2(\hat{\Gamma}_\mu^\infty(M))$ is also a Borel set. Note that for any $v \in \hat{\Gamma}_\mu^\infty(M)$ we have $L(v) = \infty$, which means that all \tilde{v} with the same marginals as v are also in $\hat{\Gamma}_\mu^\infty(M)$. This implies $\hat{\Gamma}_\mu^\infty(M) = f_2^{-1}(f_2(\hat{\Gamma}_\mu^\infty(M)))$. Hence, $\hat{\Gamma}_\mu^\infty(M)$ is a Borel set. \square

II) Proof under Assumption 2

It is easy to check that the proof under Assumption 1 remains valid if X and Y are arbitrary uncountable Polish spaces such that Y is compact, and the distortion measure ρ is an extended real valued function (no steps exploited the special structure of \mathbb{R}^n). Let Y be the one-point compactification of \mathbb{R}^n [21]. Y is clearly an uncountable Polish space. Define the extended real valued distortion measure $\rho : \mathsf{X} \times \mathsf{Y} \rightarrow [0, \infty]$ by

$$\rho(x, y) = \begin{cases} \|x - y\|^2, & \text{if } y \in \mathbb{R}^n \\ \infty, & \text{if } y = \infty. \end{cases} \quad (19)$$

It is straightforward to check that ρ is continuous. Define L on $\mathcal{P}(\mathsf{X} \times \mathsf{Y})$ as before, but with this new distortion measure ρ . The proof under Assumption 1 gives a minimizer $v^* = \int_{\Gamma_\mu(M)} v P(dv)$ for **(P1)**. Define $\hat{\Gamma}_\mu(M) = \{v \in \Gamma_\mu(M) : v(\mathsf{X} \times \{\infty\}) = 0\}$. Since $L(v^*) < \infty$ by assumption, $P(\hat{\Gamma}_\mu(M)) = 1$. This implies that v^* is also a minimizer for the problem **(P1)** when $\mathsf{X} = \mathsf{Y} = \mathbb{R}^n$ and $\rho = \|x - y\|^2$. \square

E. Proof of Theorem 4

From the proof of Theorem 3 recall the set $\hat{\Gamma}_\mu(M)$ of probability measures which solve the optimal mass transport problem for fixed input marginal μ and some output marginal ψ_0 in $\mathcal{P}_M(\mathcal{Y})$. It is known that if μ admits a density and $\rho(x, y) = \|x - y\|^2$, then each $v \in \hat{\Gamma}_\mu(M)$ is in the form $v(dx dy) = \mu(dx) \delta_{q(x)}(dy)$ for some $q \in \mathcal{Q}_{M,c}$ (see, e.g. [38, Theorem 1]). Thus in this case $\hat{\Gamma}_\mu(M) \subset \Gamma_\mu(M)$, which implies that $\hat{\Gamma}_{\mu,\psi}^{\text{R}}(M) \subset \Gamma_{\mu,\psi}^{\text{R},c}(M) \subset \Gamma_{\mu,\psi}^{\text{R}}(M)$. Recall the problem **(P2)** in the proof of Theorem 3. It was shown that **(P2)** has a minimizer v^* . It is clear from the previous discussion that v^* is obtained by randomizing over the set of quantizers having convex codecells represented by $\hat{\Gamma}_\mu(M)$. On the other hand, v^* is also a minimizer for the problem **(P1)** by Proposition 2 in the proof of Theorem 3. \square

F. Proof of Theorem 5

Recall the continuous mapping $h : \mathcal{P}(\Gamma_\mu(M)) \rightarrow \Gamma_\mu^{\text{R}}(M)$ defined in (13). Let $\mathcal{P}_F(\Gamma_\mu(M))$ denote the set of probability measures on $\Gamma_\mu(M)$ having finite support. Clearly $h(\mathcal{P}_F(\Gamma_\mu(M))) = \Gamma_\mu^{\text{FR}}(M)$.

Lemma 12. $\Gamma_\mu^{\text{FR}}(M)$ is dense in $\Gamma_\mu^{\text{R}}(M)$.

Proof: Since $\Gamma_\mu(M)$ is a separable metric space, $\mathcal{P}_F(\Gamma_\mu(M))$ is dense in $\mathcal{P}(\Gamma_\mu(M))$ by [34, Theorem 6.3]. Since $\Gamma_\mu^{\text{FR}}(M)$ is the image of a $\mathcal{P}_F(\Gamma_\mu(M))$ under the continuous function h which maps $\mathcal{P}(\Gamma_\mu(M))$ onto $\Gamma_\mu^{\text{R}}(M)$, it is dense in $\Gamma_\mu^{\text{R}}(M)$. \square

Recall that the Prokhorov metric on $\mathcal{P}(\mathsf{E})$, where (E, d) is a metric space, is defined as [14]

$$d_P(v, \nu) = \inf \left\{ \alpha : v(A) \leq \nu(A^\alpha) + \alpha, \right. \\ \left. \nu(A) \leq \nu(A^\alpha) + \alpha \text{ for all } A \in \mathcal{B}(\mathsf{E}) \right\} \quad (20)$$

where

$$A^\alpha = \left\{ e \in \mathsf{E} : \inf_{e' \in A} d(e, e') < \alpha \right\}.$$

Hence for $v, \nu \in \mathcal{P}(\mathsf{X} \times \mathsf{Y})$,

$$d_P(v, \nu) \geq \inf \left\{ \alpha : v(\mathsf{X} \times B) \leq \nu((\mathsf{X} \times B)^\alpha) + \alpha, \right. \\ \left. \nu(\mathsf{X} \times B) \leq \nu((\mathsf{X} \times B)^\alpha) + \alpha, B \in \mathcal{B}(\mathsf{Y}) \right\} \\ = d_P(v(\mathsf{X} \times \cdot), \nu(\mathsf{X} \times \cdot))$$

(note that $(\mathsf{X} \times B)^\alpha = \mathsf{X} \times B^\alpha$). This implies

$$G_\psi^\alpha := \left\{ v \in \mathcal{P}(\mathsf{X} \times \mathsf{Y}) : v(\mathsf{X} \times \cdot) \in B(\psi, \alpha) \right\} \\ \supset \left\{ v \in \mathcal{P}(\mathsf{X} \times \mathsf{Y}) : d_P(\hat{v}, v) < \alpha \right\}, \quad (21)$$

where \hat{v} is such that $\hat{v}(\mathsf{X} \times \cdot) = \psi$ and $\alpha > 0$. Recall that given a metric space E and $A \subset \mathsf{E}$, a set $B \subset A$ is relatively open in A if $B = A \cap U$ for some open set $U \subset \mathsf{E}$.

Lemma 13. $\mathcal{M}_{\mu,\psi}^\delta$ is relatively open in $\Gamma_\mu^{\text{R}}(M)$.

Proof: Since $\mathcal{M}_{\mu,\psi}^\delta = G_\psi^\delta \cap \Gamma_\mu^{\text{R}}(M)$, it is enough to prove that G_ψ^δ is open in $\mathcal{P}(\mathsf{X} \times \mathsf{Y})$. Let $\tilde{v} \in G_\psi^\delta$. Then $\tilde{v}(\mathsf{X} \times \cdot) \in B(\psi, \delta)$ by definition, and there exists $\delta_0 > 0$ such that $B(\tilde{v}(\mathsf{X} \times \cdot), \delta_0) \subset B(\psi, \delta)$. By (21) we have

$$\left\{ v \in \mathcal{P}(\mathsf{X} \times \mathsf{Y}) : d_P(\tilde{v}, v) < \delta_0 \right\} \subset G_{v(\mathsf{X} \times \cdot)}^{\delta_0}. \quad (22)$$

We also have $G_{v(\mathsf{X} \times \cdot)}^{\delta_0} \subset G_\psi^\delta$ since $B(\tilde{v}(\mathsf{X} \times \cdot), \delta_0) \subset B(\psi, \delta)$. This implies that G_ψ^δ is open in $\mathcal{P}(\mathsf{X} \times \mathsf{Y})$. \square

I) Case 1

First we treat the case $L(v) > \inf_{v' \in \Gamma_\mu(M)} L(v')$. If ρ is continuous and bounded, then L is continuous. Hence, $\{v' \in \Gamma_\mu^{\text{R}}(M) : L(v') < L(v)\}$ is relatively open in $\Gamma_\mu^{\text{R}}(M)$. Define $F := \{v' \in \Gamma_\mu^{\text{R}}(M) : L(v') < L(v)\}$.

Lemma 14. $F \cap \mathcal{M}_{\mu,\psi}^\delta$ is nonempty and relatively open in $\Gamma_\mu^{\text{R}}(M)$.

Proof: By Lemma 13 and the above discussion the intersection is clearly relatively open in $\Gamma_\mu^{\text{R}}(M)$, so we need to show that it is not empty. Since $L(v) > \inf_{v' \in \Gamma_\mu(M)} L(v')$, there exists $\tilde{v} \in \Gamma_\mu(M)$ such that $L(\tilde{v}) < L(v)$. Define the sequence of randomized quantizers $\{v_n\} \in \Gamma_\mu^{\text{R}}(M)$ by letting $v_n = \frac{1}{n} \tilde{v} + (1 - \frac{1}{n})v$. Then, $v_n \rightarrow v$ weakly because for any continuous and bounded real function f on $\mathsf{X} \times \mathsf{Y}$

$$\lim_{n \rightarrow \infty} \left| \int_{\mathsf{X} \times \mathsf{Y}} f dv_n - \int_{\mathsf{X} \times \mathsf{Y}} f dv \right| \\ = \lim_{n \rightarrow \infty} \frac{1}{n} \left| \int_{\mathsf{X} \times \mathsf{Y}} f d\tilde{v} - \int_{\mathsf{X} \times \mathsf{Y}} f dv \right| = 0.$$

Hence there exists n_0 such that $v_n \in M_{\mu,\psi}^\delta$ for all $n \geq n_0$. On the other hand, for any n

$$\begin{aligned} L(v_n) &= L\left(\frac{1}{n}\tilde{v} + \left(1 - \frac{1}{n}\right)v\right) \\ &= \frac{1}{n}L(\tilde{v}) + \left(1 - \frac{1}{n}\right)L(v) \\ &< L(v). \end{aligned}$$

This implies $v_n \in \mathcal{M}_{\mu,\psi}^\delta \cap F$ for all $n \geq n_0$, completing the proof. \square

Hence, we can conclude that there exists finitely randomized quantizer $v_F \in F \cap M_{\mu,\psi}^\delta$ by Lemmas 12 and 14. By the definition of F we also have $L(v_F) < L(v)$. This completes the proof of the theorem for this case.

II) Case 2

The case $L(v) = \inf_{v' \in \Gamma_\mu(M)} L(v') := L^*$ is handled similarly. Define the subset of $\Gamma_\mu(M)$ whose elements correspond to optimal quantizers:

$$\Gamma_{\mu,\text{opt}}(M) = \{v' \in \Gamma_\mu(M) : L(v') = L^*\}.$$

Define $\Gamma_{\mu,\text{opt}}(M) = L^{-1}(L^*) \cap \Gamma_\mu(M)$ and let $\Gamma_{\mu,\text{opt}}^R(M)$ be the randomization of $\Gamma_{\mu,\text{opt}}(M)$, obtained by replacing $\Gamma_\mu(M)$ with $\Gamma_{\mu,\text{opt}}(M)$ in (4). Note that if $L(v) = L^*$, then v is obtained by randomizing over the set $\Gamma_{\mu,\text{opt}}(M)$, i.e., $v \in \Gamma_{\mu,\text{opt}}^R(M)$. Let $\Gamma_{\mu,\text{opt}}^{\text{FR}}(M)$ denote the set obtained by the finite randomization of $\Gamma_{\mu,\text{opt}}(M)$. By using the same proof method as in Lemma 12 we can prove that $\Gamma_{\mu,\text{opt}}^{\text{FR}}(M)$ is dense in $\Gamma_{\mu,\text{opt}}^R(M)$. In addition, $\mathcal{M}_{\mu,\psi}^\delta$ is relatively open in $\Gamma_{\mu,\text{opt}}^R(M)$ by Lemma 13. Thus, there exists finitely randomized quantizer $v_F \in \mathcal{M}_{\mu,\psi}^\delta \cap \Gamma_{\mu,\text{opt}}^R(M)$ with $L(v_F) = L(v) = L^*$. This completes the proof of Theorem 5. \square

G. Proof of Theorem 6

Let $\hat{v} \in \mathcal{M}_{\mu,\psi}^\delta$ be such that $L(\hat{v}) < \inf_{v \in \mathcal{M}_{\mu,\psi}^\delta} L(v) + \varepsilon/2$. Let \hat{P} be the probability measure on $\Gamma_\mu(M)$ that induces \hat{v} , i.e., $\hat{v} = \int_{\Gamma_\mu(M)} v \hat{P}(dv)$. Consider a sequence of independent and identically distributed (i.i.d.) random variables $X_1, X_1, \dots, X_n, \dots$ defined on some probability space $(\Omega, \mathcal{F}, \gamma)$ which take values in $(\Gamma_\mu(M), \mathcal{B}(\Gamma_\mu(M)))$ and have common distribution \hat{P} . Then $L(X_1), L(X_2), \dots$ are i.i.d. \mathbb{R} -valued random variables with distribution $\hat{P} \circ L^{-1}$. Thus we have

$$\begin{aligned} \int_{\Omega} L(X_i(\omega)) \gamma(d\omega) &= \int_{\Gamma_\mu(M)} L(v) \hat{P}(dv) = L(\hat{v}) \\ &< \inf_{v \in \mathcal{M}_{\mu,\psi}^\delta} L(v) + \frac{\varepsilon}{2} \end{aligned}$$

by assumption. The empirical measures P_n^ω on $\Gamma_\mu(M)$ corresponding to X_1, \dots, X_n are

$$P_n^\omega(\cdot) := \frac{1}{n} \sum_{i=1}^n \delta_{X_i(\omega)}(\cdot).$$

By the strong law of large numbers

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n L(X_i) &= \int_{\Gamma_\mu(M)} L(v) P_n^\omega(dv) \\ &\rightarrow \int_{\Gamma_\mu(M)} L(v) \hat{P}(dv) = L(\hat{v}) \end{aligned} \quad (23)$$

γ -a.s. As a subset of $\mathcal{P}(X \times Y)$, $\Gamma_\mu(M)$ with the Prokhorov metric is a separable metric space, and thus by [21, Theorem 11.4.1] we also have the almost sure convergence of empirical measures, i.e., $P_n^\omega \rightarrow \hat{P}$ weakly γ -a.s. Thus there exists $\hat{\omega} \in \Omega$ for which both convergence results hold. Define the sequence of finitely randomized quantizers $\{v_n\}$ by $v_n = \int_{\Gamma_\mu(M)} v P_n^{\hat{\omega}}(dv)$. By (23) $L(v_n) \rightarrow L(\hat{v})$ and by Lemma 10 in the proof of Theorem 3 $v_n \rightarrow \hat{v}$ weakly. Since $\mathcal{M}_{\mu,\psi}^\delta$ is a relatively open neighborhood of \hat{v} in $\Gamma_\mu^R(M)$, we can find sufficiently large n such that $v_n \in \mathcal{M}_{\mu,\psi}^\delta$ and $L(v_n) < L(\hat{v}) + \frac{\varepsilon}{2}$. Hence, for any $\varepsilon > 0$ there exists an ε -optimal finitely randomized quantizer for **(P3)**. \square

H. Proofs for Section V

Proof of Lemma 2: The proof uses standard notation for information quantities [24]. Let $X^n \sim \mu^n$, $Z \sim \nu$, and $Y^n = \mathfrak{q}(X^n, Z) \sim \psi^n$, where (\mathfrak{q}, ν) is an arbitrary Model 2 randomized quantizer with at most 2^{nR} levels (Z is independent of X^n). Let $D_i = E[\rho(X_i, Y_i)]$ and $D = \frac{1}{n} \sum_{i=1}^n D_i = E[\rho_n(X^n, Y^n)]$. Since $\mathfrak{q}(\cdot, z)$ has at most 2^{nR} levels for each z ,

$$\begin{aligned} nR &\geq H(Y^n|Z) \geq I(X^n; Y^n|Z) \\ &\geq I(X^n; Y^n) \end{aligned} \quad (24)$$

$$\begin{aligned} &\geq \sum_{i=1}^n I(X_i; Y_i) \\ &\geq \sum_{i=1}^n I_m(\mu\|\psi, D_i) \\ &\geq nI_m(\mu\|\psi, D) \end{aligned} \quad (25)$$

where in the last two inequalities follow since $Y_i \sim \psi$, $i = 1, \dots, n$ and $I_m(\mu\|\psi, D)$ is convex in D [23, Appendix A]. Inequalities (24) and (25) follow from the chain rule for mutual information (Kolmogorov's formula) [39, Corollary 7.14], which in particular implies that $I(U; V|W) \geq I(U; V)$ for general random variables U, V , and W , defined on the same probability space, such that U and W are independent. This proves that $R \geq I_m(\mu\|\psi, D)$. \square

Proof of Lemma 3: Let $U^{2^{nR}} = (U^n(1), \dots, U^n(2^{nR}))$ which is a $n2^{nR}$ -vector. Then, we can write

$$\hat{X}^n = g(X^n, U^{2^{nR}})$$

for a function g from $Y^{n(2^{nR}+1)}$ to Y^n . Observe the following:

- (i) For any permutation σ of $\{1, \dots, n\}$, X^n and $X_\sigma^n = (X_{\sigma(1)}, \dots, X_{\sigma(n)})$ have the same distribution. The same issue is true for $U^n(i)$ and $U^n(i)_\sigma$ for all i because for any $u^n \in T_n(\psi_n)$, $u_\sigma^n \in T_n(\psi_n)$ and this mapping is a bijection on $T_n(\psi_n)$. It follows from the independence of X^n and $U^n(i)$ that (X^n, U^{nR}) and $(X_\sigma^n, U_\sigma^{2^{nR}})$ have the same distribution, where $U_\sigma^{2^{nR}} := (U_\sigma^n(1)_\sigma, \dots, U_\sigma^n(2^{nR})_\sigma)$. Thus, $g(X^n, U^{2^{nR}})$ and $g(X_\sigma^n, U_\sigma^{2^{nR}})$ have the same distribution.
- (ii) For any $x^n \in X^n$ and $y^n \in Y^n$, $\rho_n(x^n, y^n) = \rho_n(x_\sigma^n, y_\sigma^n)$. Thus, if g outputs $u^n(i)$ for inputs

$x^n, u^n(1), \dots, u^n(2^{nR})$, then g outputs $u^n(i)_\sigma$ for inputs $x_\sigma^n, u^n(1)_\sigma, \dots, u^n(2^{nR})_\sigma$. It follows that

$$g(X_\sigma^n, U_\sigma^{2^{nR}}) = g(X^n, U^{2^{nR}})_\sigma.$$

Together with $i)$ this implies that \hat{X}^n and \hat{X}_σ^n have the same distribution.

Let u^n and $v^n \in T_n(\psi_n^{(n)})$ and so $u^n = v_\sigma^n$ for some permutation σ . Then (ii) implies

$$\Pr\{\hat{X}^n = u^n\} = \Pr\{\hat{X}_\sigma^n = u^n\}.$$

Since $\Pr\{\hat{X}^n = v^n\} = \Pr\{\hat{X}_\sigma^n = v_\sigma^n\}$ and $v_\sigma^n = u^n$, we obtain

$$\Pr\{\hat{X}^n = u^n\} = \Pr\{\hat{X}^n = v^n\}$$

proving that \hat{X}^n is uniform on $T_n(\psi_n^{(n)})$. \square

Proof of Lemma 4: By [24, Theorem 11.1.2] we have

$$\begin{aligned} \frac{1}{n} \mathcal{D}(\psi_n^{(n)} \parallel \psi^n) &= \frac{1}{n} \sum_{y^n \in T_n(\psi_n)} \psi_n^{(n)}(y^n) \log \frac{\psi_n^{(n)}(y^n)}{\psi^n(y^n)} \\ &= \frac{1}{n} \log \frac{2^{n(H(\psi_n) + \mathcal{D}(\psi_n \parallel \psi))}}{|T_n(\psi_n)|}. \end{aligned} \quad (26)$$

From [24, Theorem 11.1.3],

$$\frac{1}{(n+1)^{|\mathcal{X}|}} 2^{nH(\psi_n)} \leq |T_n(\psi_n)| \leq 2^{nH(\psi_n)}$$

and thus $\frac{1}{n} \mathcal{D}(\psi_n^{(n)} \parallel \psi^n)$ is sandwiched between $\mathcal{D}(\psi_n \parallel \psi)$ and $\frac{|\mathcal{X}|}{n} \log(n+1) + \mathcal{D}(\psi_n \parallel \psi)$. Thus

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathcal{D}(\psi_n^{(n)} \parallel \psi^n) = \lim_{n \rightarrow \infty} \mathcal{D}(\psi_n \parallel \psi) = 0$$

where the second limit holds since \mathcal{X} is a finite set and $\psi_n \rightarrow \psi$ in the l_1 -distance. \square

Proof of Lemma 5:

Let ρ^H denote the Hamming distortion and let $\rho_n^H(x^n, y^n) = (1/n) \sum_{i=1}^n \rho^H(x_i, y_i)$. Since $\rho(x, x) = 0$ for all $x \in \mathcal{X}$, we have

$$\rho_n(x^n, y^n) \leq \rho_{\max} \rho_n^H(x^n, y^n).$$

Let $T_n^H(\psi_n^{(n)}, \psi^n)$ be the distortion of the optimal coupling between $\psi_n^{(n)}$ and ψ^n when the cost function is ρ_n^H . Then the above inequality gives

$$\hat{T}_n(\psi_n^{(n)}, \psi^n) \leq \rho_{\max} T_n^H(\psi_n^{(n)}, \psi^n).$$

On the other hand, by Marton's inequality [25, Proposition 1]

$$T_n^H(\psi_n^{(n)}, \psi^n) \leq \sqrt{\frac{1}{2n} \mathcal{D}(\psi_n^{(n)} \parallel \psi^n)}.$$

Combining these bounds with $\frac{1}{n} \mathcal{D}(\psi_n^{(n)} \parallel \psi^n) \rightarrow 0$ (Lemma 4), we obtain

$$\lim_{n \rightarrow \infty} \hat{T}_n(\psi_n^{(n)}, \psi^n) = 0 \quad (27)$$

which is the first statement of the lemma.

Recall that $\rho(x, y) = d(x, y)^p$ for some $p > 0$, where d is a metric. Let $q = \max\{1, p\}$. If $p \geq 1$, then $\|V^n\|_p := (E[\sum_{i=1}^n |V_i|^p])^{1/q}$ is a norm on \mathbb{R}^n -valued random vectors whose components have finite p th moments, and if $1 < p < 0$,

we still have $\|U^n + V^n\|_p \leq \|U^n\|_p + \|V^n\|_p$. Thus we can upper bound $E[\rho_n(X^n, Y^n)]$ as follows:

$$\begin{aligned} &\left(E \left[\frac{1}{n} \sum_{i=1}^n \rho(X_i, Y_i) \right] \right)^{1/q} \\ &= \left(E \left[\frac{1}{n} \sum_{i=1}^n d(X_i, Y_i)^p \right] \right)^{1/q} \\ &\leq \left(E \left[\frac{1}{n} \sum_{i=1}^n d(X_i, \hat{X}_i)^p \right] \right)^{1/q} + \left(E \left[\frac{1}{n} \sum_{i=1}^n d(\hat{X}_i, Y_i)^p \right] \right)^{1/q} \\ &= \left(E[\rho_n(X^n, \hat{X}^n)] \right)^{1/q} + \hat{T}_n(\psi_n^{(n)}, \psi^n)^{1/q}. \end{aligned}$$

Hence (8) and (27) imply

$$\limsup_{n \rightarrow \infty} E[\rho_n(X^n, Y^n)] \leq D$$

as claimed. \square

Proof of Lemma 6: Let $X \sim \mu$ and $Y \sim \psi$ such that $I(X; Y)$ achieves $I_m(\mu \parallel \psi, D) < \infty$ at distortion level D (the existence of such pair follows from an analogous statements for rate-distortion functions [40]). Let q_k denote the uniform quantizer on the interval $[-k, k]$ having 2^k levels, where we extend q_k to the real line by using the nearest neighborhood encoding rule. Let $X(k) = q_k(X)$ and $Y(k) = q_k(Y)$. We clearly have

$$E[(X - X(k))^2] \rightarrow 0, \quad E[(Y - Y(k))^2] \rightarrow 0 \text{ as } k \rightarrow \infty. \quad (28)$$

Let μ_k and ψ_k denote the distributions of $X(k)$ and $Y(k)$, respectively. Then by [22, Theorem 6.9] it follows that $\hat{T}_1(\mu_k, \mu) \rightarrow 0$ and $\hat{T}_1(\psi_k, \psi) \rightarrow 0$ as $k \rightarrow \infty$ since $\mu_k \rightarrow \mu$, $\psi_k \rightarrow \psi$ weakly, and $E[X(k)^2] \rightarrow E[X^2]$, $E[Y(k)^2] \rightarrow E[Y^2]$.

By the data processing inequality, we have for all k ,

$$I(X(k); Y(k)) \leq I(X; Y). \quad (29)$$

Also note that (28) implies

$$\begin{aligned} &\limsup_{k \rightarrow \infty} E[\rho_1(X(k), Y(k))] \\ &= \limsup_{k \rightarrow \infty} E[(X(k) - Y(k))^2] \leq D. \end{aligned}$$

Thus, for given $\varepsilon > 0$, if k is large we have $I_m(\mu_k \parallel \psi_k, D + \varepsilon) \leq I_m(\mu \parallel \psi, D)$ as claimed. \square

ACKNOWLEDGEMENT

The authors would like to thank Ram Zamir for helpful discussions concerning the proof of Theorem 7. Naci Saldi would like to thank Marcos Vasconcelos for discussions on optimal transport theory.

REFERENCES

- [1] L. Roberts, "Picture coding using pseudo-random noise," *IEEE Trans. Inf. Theory*, vol. 8, no. 2, pp. 145–154, Feb. 1962.
- [2] L. Schucman, "Dither signals and their effect on quantization noise," *IEEE Trans. Commun.*, vol. 12, no. 4, pp. 162–165, Dec. 1964.
- [3] R. Gray and T. Stockham, "Dithered quantizers," *IEEE Trans. Inf. Theory*, vol. 39, no. 3, pp. 805–812, May 1993.

- [4] J. Ziv, "On universal quantization," *IEEE Trans. Inf. Theory*, vol. 31, no. 3, pp. 344–347, May 1985.
- [5] R. Zamir and M. Feder, "On universal quantization by randomized uniform/lattice quantizers," *IEEE Trans. Inf. Theory*, vol. 38, no. 2, pp. 428–436, Mar. 1992.
- [6] —, "Information rates of pre/post-filtered dithered quantizers," *IEEE Trans. Inf. Theory*, vol. 42, no. 5, pp. 1340–1353, Sep. 1996.
- [7] E. Akyol and K. Rose, "On constrained randomized quantization," in *Proc. Data Compress. Conf.*, Snowbird, Utah, USA, Apr. 2012, pp. 72–81.
- [8] —, "On constrained randomized quantization," *IEEE Trans. Signal Processing*, vol. 61, no. 13, pp. 3291–3302, Jul. 2013.
- [9] C. E. Shannon, "Coding theorems for a discrete source with a fidelity criterion," *IRE Nat. Conv. Rec.*, vol. part 4, pp. 138–143, 1959.
- [10] M. Li, J. Klejsa, and W. Kleijn, "Distribution preserving quantization with dithering and transformation," *IEEE Signal Processing Letters*, vol. 17, no. 12, pp. 1014–1017, 2010.
- [11] —, "On distribution preserving quantization," *arXiv:1108.3728*, 2011.
- [12] J. Klejsa, G. Zhang, and M. L. and W.B. Kleijn, "Multiple description distribution preserving quantization," *IEEE Trans. Signal Processing*, vol. 61, no. 24, pp. 6410–6422, Dec. 2013.
- [13] V. Borkar, "White-noise representations in stochastic realization theory," *SIAM J. Control Optim.*, vol. 31, no. 5, pp. 1093–1102, 1993.
- [14] P. Billingsley, *Convergence of probability measures*, 2nd ed. New York: Wiley, 1999.
- [15] R. Zamir, *Lattice Coding for Signals and Networks*. Oxford University Press, 2014.
- [16] V. Borkar, S. Mitter, A. Sahai, and S. Tatikonda, "Sequential source coding: An optimization viewpoint," in *Proc. IEEE Conference on Decision and Control*, Seville, Spain, Dec. 2005, pp. 1035–1042.
- [17] S. Yüksel and T. Linder, "Optimization and convergence of observation channels in stochastic control," *SIAM J. Control Optim.*, vol. 50, no. 2, pp. 864–887, 2012.
- [18] W. Kreitmeier, "Optimal vector quantization in terms of Wasserstein distance," *J. Multivariate Anal.*, vol. 102, no. 8, pp. 1225–1239, 2011.
- [19] S. Graf and H. Luschgy, *Foundations of Quantization for Probability Distributions*. Springer, 2000.
- [20] O. Hernández-Lerma and J. Lasserre, *Discrete-Time Markov Control Processes: Basic Optimality Criteria*. Springer, 1996.
- [21] R. M. Dudley, *Real Analysis and Probability*. New York: Chapman and Hall, 1989.
- [22] C. Villani, *Optimal transport: old and new*. Springer, 2009.
- [23] R. Zamir and K. Rose, "Natural type selection in adaptive lossy compression," *IEEE Trans. Inf. Theory*, vol. 47, no. 1, pp. 99–111, Jan. 2001.
- [24] T. Cover and J. Thomas, *Elements of Information Theory*, 2nd ed. Wiley, 2006.
- [25] K. Marton, "Bounding \bar{d} -distance by informational divergence: a method to prove measure concentration," *Ann. Probab.*, vol. 24, no. 2, pp. 857–866, 1996.
- [26] E.-H. Yang and J. C. Kieffer, "On the performance of data compression algorithms based upon string matching," *IEEE Trans. Inform. Theory*, vol. 44, no. 1, pp. 47–65, Jan. 1998.
- [27] T. Weissman and E. Ordentlich, "The empirical distribution of rate-constrained source codes," *IEEE Trans. Inform. Theory*, vol. 51, no. 11, pp. 3718–3733, Nov. 2005.
- [28] V. P. Crawford and J. Sobel, "Strategic information transmission," *Econometrica*, vol. 50, pp. 1431–1451, 1982.
- [29] T. Linder and S. Yüksel, "On optimal zero-delay coding of vector Markov sources," *IEEE Trans. Inform. Theory*, vol. 60, no. 10, pp. 5975–5991, Oct. 2014.
- [30] S. Yüksel and T. Başar, *Stochastic Networked Control Systems: Stabilization and Optimization under Information Constraints*. Boston, MA: Birkhäuser, 2013.
- [31] D. P. Bertsekas and S. E. Shreve, *Stochastic optimal control: The discrete time case*. Academic Press New York, 1978.
- [32] L. Dubins and D. Freedman, "Measurable sets of measures," *Pacific J. Math*, vol. 14, no. 4, pp. 1211–1222, 1964.
- [33] P. Billingsley, *Probability and Measure*, 3rd ed. Wiley, 1995.
- [34] K. Parthasarathy, *Probability Measures on Metric Spaces*. AMS Bookstore, 1967.
- [35] E. B. Dykin, *Controlled Markov Processes*. Berlin, New York: Springer-Verlag, 1979.
- [36] V. Borkar, "On extremal solutions to stochastic control problems," *Appl. Math. Optim.*, vol. 24, no. 1, pp. 317–330, 1991.
- [37] A. Pratelli, "On the sufficiency of c -cyclical monotonicity for optimality of transport plans," *Math. Z.*, vol. 258, no. 3, pp. 677–690, 2008.
- [38] M. McAsey and L. Mou, "Optimal locations and the mass transport problem," *Contemp. Math.*, vol. 226, pp. 131–148, 1999.
- [39] R. Gray, *Entropy and Information Theory*. Springer, 2011.
- [40] I. Csiszár, "On an extremum problem of information theory," *Studia Scientiarum Mathematicarum Hungarica*, pp. 57–70, 1974.