

Price of metric universality in vector quantization is at most 0.11 bit

Alina Harbuzova

Massachusetts Institute of Technology. Supported by MathWorks Fellowship and Siebel Scholarship

Or Ordentlich

Hebrew University of Jerusalem. Supported by the Israel Science Foundation (ISF), grant No. 2878/25

Yury Polyanskiy

Massachusetts Institute of Technology. Supported by NSF Grant No. 2112665 via subaward KR 704702 from the University of California, San Diego

Editors: Steve Hanneke and Tor Lattimore

Abstract

Fast computation of a matrix product $W^\top X$ is a workhorse of modern LLMs. To make their deployment more efficient, a popular approach is that of using a low-precision approximation \widehat{W} in place of true W (“weight-only quantization”). Information theory demonstrates that an optimal algorithm for reducing precision of W depends on the (second order) statistics of X and requires a careful alignment of vector quantization codebook with PCA directions of X (a process known as “waterfilling allocation”). Dependence of the codebook on statistics of X , however, is highly impractical. This paper proves that there exist a universal codebook that is simultaneously near-optimal for all possible statistics of X , in the sense of being at least as good as an X -adapted waterfilling codebook with rate reduced by 0.11 bit per dimension in the case when W is Gaussian. Such universal codebook would be an ideal candidate for the low-precision storage format, a topic of active modern research, but alas the existence proof is non-constructive.

Equivalently, our result shows existence of a net in \mathbb{R}^n that is a nearly-optimal covering of a sphere simultaneously with respect to all Hilbert norms.

Keywords: vector quantization, oracle bounds, rate-distortion, waterfilling, regret, universality

1. Introduction

The most basic element of all modern AI is a neural unit: given a (dynamically changing) activation vector $X \in \mathbb{R}^n$ the unit needs to compute the output

$$Y = W^\top X,$$

where $W \in \mathbb{R}^n$ is a (static) weight. The problem, rapidly becoming central for economic deployment and continued evolution of large language models (LLMs), is to reduce storage/communication requirement by saving W in “low-precision”. (The symmetric question of also converting X to low-precision is outside of scope of this work, though see [Ordentlich and Polyanskiy \(2025\)](#) for some recent theoretical analysis.)

Early deep learning models saved W in full-precision (known as FP32, and corresponding to $R = 32$ bit / coordinate), but soon moved to half-precision (FP16/BF16, or $R = 16$ bit). In the domain of LLMs, pioneering work [Dettmers et al. \(2022\)](#) showed that very little degradation is introduced if W is approximated by rescaling it to appropriate range and then rounding each coordinate of normalized W to nearest integer in $\{-128, -127, \dots, 127\}$, a so called *INT8* quantization.

Subsequently, more sophisticated low-precision storage formats were introduced, with currently the most popular being NVFP4 (rate $R = 4.5$ bit) and MXFP4 (rate $R = 4.25$ bit), see [NVIDIA et al. \(2025\)](#); [Open Compute Project \(2023\)](#).

In this paper we are focusing on a fundamental question: what is the best way of reducing precision of W ? That is, how to replace W by a version \widehat{W} that incurs minimal degradation of performance, i.e. $\widehat{W}^\top X \approx W^\top X$, while admitting a short bit-length description. A natural way to do that, known as *vector quantization*, would be to pre-define a *codebook* $\mathbf{C} \subset \mathbb{R}^n$ of size $|\mathbf{C}| = 2^{nR}$. Clearly, any element of \mathbf{C} can be described by nR bits, hence achieving rate R of bits / coordinate. Given \mathbf{C} we approximate W as

$$\widehat{W} = \underset{c \in \mathbf{C}}{\operatorname{argmin}} d(W, c)$$

for some distance metric $d(\cdot, \cdot)$. When d is a standard Euclidean metric, then the problem reduces to a classical vector quantization problem in \mathbb{R}^n , with many classical solutions including lattices and trellis-coded constructions [Gersho and Gray \(2012\)](#). However, as was brilliantly shown by [Frantar et al. \(2023\)](#) large savings can be made if metric $d(\cdot, \cdot)$ is chosen with the knowledge of statistics of X in mind.

Indeed, if X is modeled as random with second-order statistics $\Sigma_X = \mathbb{E}[XX^\top]$ then

$$\mathbb{E}_X \left[(Y - \widehat{Y})^2 \right] = \mathbb{E}_X \left[(W^\top X - \widehat{W}^\top X)^2 \right] = (W - \widehat{W})^\top \Sigma_X (W - \widehat{W}).$$

Thus, we see that a natural choice of metric (given knowledge of Σ_X) is

$$d_{\Sigma_X}(W, \widehat{W}) = \mathbb{E}_X \left[(W^\top X - \widehat{W}^\top X)^2 \right] = (W - \widehat{W})^\top \Sigma_X (W - \widehat{W}). \quad (d_{\Sigma_X})$$

The easiest way to demonstrate how adaptation to Σ_X can significantly improve rate-distortion tradeoff is to consider a rank-1 case, i.e. when X is always collinear with a fixed vector $v \in \mathbb{R}^n$. In this case, a clever choice of the codebook \mathbf{C} is $\{0, \pm \epsilon v, \pm 2\epsilon v, \dots\}$, i.e. very fine quantization along a single direction v in \mathbb{R}^n . Indeed, by not needing to spread the points of \mathbf{C} among all n dimensions, one can get exponential improvement in quality of approximation of $W^\top X$, since only the scalar value $W^\top v$ affects the result. Since activations in LLMs are notoriously low-rank, this adaptation of \mathbf{C} to directions of principal variation (PCA) of X understandably improves performance.

Herein, however, lies the main problem that we are trying to address: while adapting \mathbf{C} to the statistics of the input X is desirable, it may not be generally possible due to restrictions of hardware. Indeed, the mapping from actual bits (loaded from memory) to elements of \mathbf{C} needs to be fixed at hardware design stage and cannot depend on statistics of X (in particular, because the same hardware is used for implementing different neurons, facing different types of X). Below we call this requirement, alternatively, as *universal codebook* or a Σ_X -*oblivious decoder*, to reflect the fact that \mathbf{C} has to be universal across all possible choices of Σ_X .

To continue with more quantitative investigation, let us make a modeling assumption (well justified by empirical statistics of LLM weight matrices) that $W \sim \mathcal{N}(0, I_n)$. In this case a given codebook \mathbf{C} under Σ_X statistics attains *distortion*:

$$D(\mathbf{C}, \Sigma_X) := \frac{1}{n} \mathbb{E}_W \left[\min_{c \in \mathbf{C}} d_{\Sigma_X}(W, c) \right]. \quad (1)$$

Classical information-theoretic field, known as rate-distortion theory, establishes that for a codebook \mathbf{C} to achieve distortion $D(\mathbf{C}, \Sigma_X) \leq D$ one must have

$$\log |\mathbf{C}| \geq n\mathbf{R}_{\text{wf}}(\Sigma_X, D),$$

where \mathbf{R}_{wf} is given by a so-called *waterfilling formula*, see Prop. 2.

The main question of this work: *How much does the requirement of universality cost in terms of performance?* For example, in the rank-1 case above the waterfilling codebook would allocate its elements along a single direction v . This codebook, however, would be grossly suboptimal for another rank-1 Σ_X which has its PCA direction orthogonal to v . Somewhat surprisingly, thus, we show that nevertheless the answer is *not much*. The main result of this work is demonstration of existence of a universal \mathbf{C} , which is simultaneously near optimal for all possible Σ_X . Informally, we can state our main result as follows.

Theorem 1 (Informal: Universality costs ≤ 0.11 Bits) *Let us assume that $W \sim \mathcal{N}(0, I_n)$ and let $\mathbf{R}_{\text{wf}}(\Sigma_X, D)$ denote the information-theoretic (waterfilling) lower bound on rate needed to achieve distortion at most D in the oracle setting, where codebook is optimized for a fixed Σ_X . There exists a universal codebook \mathbf{C} with 2^{nR} points such that its distortion simultaneously for all $\Sigma_X \in \mathbb{S}_+^n$ satisfies:*

$$R \leq \mathbf{R}_{\text{wf}}(\Sigma_X, D(\mathbf{C}, \Sigma_X)) + 0.11 \text{ bit}.$$

The implication for hardware design is clear: it is possible to create a universal low-precision storage format (for W) that is optimal (up to rate gap of at most 0.11 bit) simultaneously for all kinds of distributions of statistics of the other factor (X) in the inner-product. Our result can also be interpreted as a statement about metric entropy: *There exists a universal net on a unit sphere, which covers unit sphere near-optimally simultaneously for all possible Hilbert norms on \mathbb{R}^n .*

Paper organization. The following Section 2 formalizes weight-only quantization for inner products under the distortion d_{Σ_X} . We then introduce the oracle benchmark given by the Gaussian rate-distortion tradeoff under weighted MSE, attained by the waterfilling solution in the setting where both encoder and decoder know the second-order statistics Σ_X (Prop. 2). Our main results are stated in Theorems 3 and 5. Theorem 3 shows the existence of a universal codebook achieving the explicit rate-distortion tradeoff (RDRC) over all Σ_X . Theorem 5 upper bounds the worst-case rate overhead incurred by this universal decoder relative to the oracle waterfilling benchmark.

Section 3 provides a proof sketch of Theorem 3 and the main geometric ideas: Section 3.1 describes the universal codebook construction and intuition, and Section 3.2 outlines the random-coding analysis leading to (RDRC).

Complete proofs are given in the Appendix. While the results in Section 2 are presented for $W \sim \mathcal{N}(0, I_n)$, we first establish a general bound for a fixed (non-random) W in Section B. This result is then specialized to $W \sim \mathcal{N}(0, I_n)$ using concentration and covering argument in Section C, completing the proof of Theorem 3. Theorem 5 is proved in Section D: we derive explicit expressions for the rate-gap and prove that the maximum gap occurs at spectra with at most 2 distinct eigenvalues and vanishing distortions.

2. Main results and discussion

Consider an arbitrary $\Sigma_X \succeq 0$ and define a (square of) Hilbert metric with respect to Σ_X as in (d_{Σ_X}). We consider the problem of obtaining a low-precision (at rate R bits per coordinate) representation

\widehat{W} of a random vector $W \in \mathbb{R}^n$ with the goal of minimizing $d_{\Sigma_X}(\widehat{W}, W)$. We focus presentation of results on the the standard setting in which W is an isotropic Gaussian vector:

$$W \sim \mathcal{N}(0, I_n),$$

though the key technical results hold for general W (Section B). An (n, R) quantization scheme consists of an encoder $f : \mathbb{R}^n \rightarrow [2^{nR}]$ and decoder $g : [2^{nR}] \rightarrow \mathbb{R}^n$ and we set $\widehat{W} = g(f(W))$. The image of g is called the codebook $\mathbf{C} := \text{im } g$. The optimal encoder consists of finding a nearest to W element of \mathbf{C} , and hence we can equivalently think of a quantization scheme as completely defined by \mathbf{C} . The distortion of \mathbf{C} for a given Σ_X is denoted $D(\mathbf{C}, \Sigma_X)$, cf. (1).

Let us start with a simple case of Σ_X fixed (and hence known to both encoder f_{Σ_X} and decoder g_{Σ_X}). In this case, the optimal tradeoff between the distortion D and rate R is given by waterfilling, which we review.

Let $\Sigma_X = U\Lambda U^\top$ be the eigendecomposition with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Consider a parametric curve

$$D_{\text{wf}}(\Sigma_X, t) = \frac{1}{n} \sum_{i=1}^n \min\{\lambda_i, t\} \quad R_{\text{wf}}(\Sigma_X, t) = \frac{1}{2n} \sum_{i=1}^n \max\{0, \log(\lambda_i/t)\}, \quad (\text{WF})$$

where t is the parameter (waterfilling level). The above implicitly defines waterfilling distortion as a function of rate:

$$\mathbf{D}_{\text{wf}}(R, \Sigma_X) \triangleq D_{\text{wf}}(\Sigma_X, t_{\text{wf}}(R)), \quad \text{where } R_{\text{wf}}(\Sigma_X, t_{\text{wf}}(R)) = R. \quad (\mathbf{D}_{\text{wf}})$$

It turns out that this function indeed determines the fundamental limits in the case of oracle-knowledge of Σ_X . More exactly, we have the following.

Proposition 2 (Waterfilling) *Let $\Sigma_X \in \mathbb{S}_+^n$. For any compression scheme of rate R we have*

$$\mathbb{E}[d_{\Sigma_X}(W, g(f(W)))] \geq n\mathbf{D}_{\text{wf}}(R, \Sigma_X). \quad (2)$$

Conversely, for any $B > 0$ there exist a $c = c(B) > 0$ such that for any Σ_X there exist f and g (both depending on Σ_X) such that

$$\mathbb{E}[d_{\Sigma_X}(W, g(f(W)))] \leq n\mathbf{D}_{\text{wf}}\left(R - c\sqrt{\frac{\log n}{n}}, \Sigma_X\right) + cn^{-B} \text{tr } \Sigma_X. \quad (3)$$

Proof We give the proof of the lower bound (2) below; since the upper bound (3) is not relevant for the rest of this paper, we defer the brief of sketch of the proof to Appendix F. Though we do emphasize that the upper bound does not follow from classical theory, which concerns with separable (additive over coordinates) distortion measures, and we have to invoke more modern single-shot bounds, cf. (Polyanskiy and Wu, 2024, Chapter 25).

Let $\Sigma_X = U\Lambda U^\top$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$, and define $W' \triangleq U^\top W$, $\widehat{W}' \triangleq U^\top \widehat{W}$. Then

$$d_{\Sigma_X}(W, \widehat{W}) = (W - \widehat{W})^\top \Sigma_X (W - \widehat{W}) = (W' - \widehat{W}')^\top \Lambda (W' - \widehat{W}') = \sum_{i=1}^n \lambda_i (W'_i - \widehat{W}'_i)^2.$$

In the oracle setting (where both the encoder and decoder know Σ_X), we may equivalently compress W' , reconstruct \widehat{W}' , and output $\widehat{W} = U\widehat{W}'$. After this coordinate change, the problem becomes that of a weighted mean squared error. Since $W' \sim \mathcal{N}(0, I_n)$ has independent coordinates, a standard data-processing argument (see (Polyanskiy and Wu, 2024, Section 23.4 and Theorem 6.1)) gives:

$$nR \geq I(W'; \widehat{W}') \geq \sum_{i=1}^n I(W'_i; \widehat{W}'_i),$$

Writing $D_i \triangleq \mathbb{E} \left[(W'_i - \widehat{W}'_i)^2 \right]$, the smallest $I(W'_i; \widehat{W}'_i) = \frac{1}{2} \log \frac{1}{D_i}$ (given the value D_i) is then attained under Gaussian coupling, cf. (Polyanskiy and Wu, 2024, Section 26.1.2), which results in

$$R \geq \frac{1}{2n} \sum_{i=1}^n \log \frac{1}{D_i} \quad \text{and} \quad \frac{1}{n} \mathbb{E} \left[d_{\Sigma_X}(W, \widehat{W}) \right] = \frac{1}{n} \sum_{i=1}^n \lambda_i D_i.$$

Minimizing $\frac{1}{n} \sum_i \lambda_i D_i$ subject to the rate constraint via Lagrange multipliers gives the (reverse) waterfilling optimum, summarized by the parametric curve in (WF). \blacksquare

Now, as we discussed above, Σ_X describes distribution of activations and, practically speaking, is usually unavailable to decoder. Indeed, even if the eigenvalues Λ were known, optimal waterfilling requires knowledge of the eigenbasis U , since the decoder g_{Σ_X} essentially computes $U \cdot \widehat{U^\top W}$, where $\widehat{U^\top W}$ is the closest codeword to $U^\top W$. Communicating $U \in \mathcal{O}(n)$ to the decoder is expensive: it requires approximating $\Theta(n^2)$ real parameters, which will consume much larger than $\Theta(n)$ bits allocated for communicating W itself.¹ Thus, the tradeoff in (WF) is unattainable via a naïve “send W and Σ_X ” scheme when decoder lacks Σ_X .

To capture the limitation above, we need to assume that Σ_X is available at *encoding* time, but unavailable at *decoding* time (since the deployed dequantizer is fixed, possible even in hardware). So, formally we define a universal (n, R) quantization scheme as a pair

$$f : \mathbb{R}^n \times \mathbb{S}_+^n \rightarrow [2^{nR}], \quad g : [2^{nR}] \rightarrow \mathbb{R}^n, \quad (f, g)$$

where the encoder takes Σ_X as an input, while the decoder has no access to Σ_X and outputs $\widehat{W} = g(f(W, \Sigma_X))$. Again, the image of g is called the codebook $\mathbf{C} = \text{im } g$, which is independent of Σ_X .

The goal is to design a pair (f, g) such that $D(\mathbf{C}, \Sigma_X)$ were low simultaneously for all Σ_X . Our first main result proves existence of a universal codebook \mathbf{C} with an explicit guarantee on the achieved distortion. To define that guarantee, again let $\Sigma_X \in \mathbb{S}_+^n$ with $\text{tr}(\Sigma_X) = n$ and spectrum $\lambda = (\lambda_1, \dots, \lambda_n)$. The *random-coding rate-distortion function* is a parametric curve (with parameter $T > 0$) given by

$$D_{\text{rc}}(\lambda, T) = \frac{1}{n} \sum_{i=1}^n \frac{\lambda_i}{1 + \lambda_i T} \quad R_{\text{rc}}(\lambda, T) = \frac{1}{2n} \sum_{i=1}^n \log(1 + \lambda_i T). \quad (\text{RDRC})$$

Denote

$$\mathbf{D}_{\text{rc}}(\lambda, R) =: D_{\text{rc}}(\lambda, T_{\text{rc}}(\lambda, R)), \quad \text{where } R_{\text{rc}}(\lambda, T_{\text{rc}}(\lambda, R)) = R. \quad (\mathbf{D}_{\text{rc}})$$

1. In practice, GPU computes not a single inner product but $W^\top X$ for $W \in \mathbb{R}^{n \times a}$ being a matrix with $a \asymp n$. Thus, some of the cost of sending U amortizes over a , but still makes it a highly suboptimal choice.

Theorem 3 (Main Result I: Universal Quantization Scheme for Gaussian Input) *Fix any constants $R^*, \varepsilon, \eta, B > 0$. There exists an encoder $f : \mathbb{R}^n \times \mathbb{S}_+^n \times [0, 1] \rightarrow [2^{nR}]$, a decoder $g : [2^{nR}] \times [0, 1] \rightarrow \mathbb{R}^n$ with $R \leq R^* + \varepsilon$, and a (shared) random variable $S \in [0, 1]$ with the following property. For $W \sim \mathcal{N}(0, I_n)$ and $\Sigma_X \in \mathbb{S}_+^n$ we set $\widehat{W}(\Sigma_X) = g(f(W, \Sigma_X, S), S)$. Then, for sufficiently large $n \geq n_0 = n_0(\varepsilon, \eta, R^*, B)$, we have with probability at least $1 - \exp(-n^B)$ over $S \sim \text{Unif}[0, 1]$ that*

$$\frac{1}{n} \mathbb{E}_W \left[d_{\Sigma_X}(W, \widehat{W}(\Sigma_X)) \right] \leq \mathbf{D}_{\text{rc}}(\text{spec}(\Sigma_X), R^*) + \eta$$

simultaneously for all Σ_X with $\text{tr} \Sigma_X = n$.

Theorem 3 uses S to generate an entire codebook \mathbf{C} and shows that with high probability it achieves $\approx \mathbf{D}_{\text{rc}}(\Sigma_X, R)$ distortion. Of course, by fixing a value of S it implies *existence* of a single codebook \mathbf{C} with the same property. Formally, we have a corollary.

Corollary 4 *Under the same assumptions as in Thm 3, there exists a codebook \mathbf{C} of size $\log_2 |\mathbf{C}| \leq n(R^* + \varepsilon)$ and an encoder-decoder pair f, g (see Eq. (f, g)) with the following property. Set $\widehat{W}(\Sigma_X) = g(f(W, \Sigma_X))$. Then, simultaneously for all Σ_X we have*

$$\frac{1}{n} \mathbb{E}_W \left[d_{\Sigma_X}(W, \widehat{W}(\Sigma_X)) \right] \leq \mathbf{D}_{\text{rc}}(\text{spec}(\Sigma_X), R^*) + \eta \frac{\text{tr} \Sigma_X}{n}.$$

The description of the scheme and intuition are in Sec. 3 and the full proof in Appendix C. While the results are presented for isotropic Gaussian W , our proof proceeds by showing a more general result for a fixed non-random W (see Sec. B). Then, the result for Gaussian W is obtained by using concentration of measure (Sec. C).

What Theorem 3 shows is that, roughly speaking, there exists a universal codebook of rate R which achieves distortion $\mathbf{D}_{\text{rc}}(\text{spec}(\Sigma_X), R)$ simultaneously for all Σ_X . A natural question is: *how far is $\mathbf{D}_{\text{rc}}(\text{spec}(\Sigma_X), R)$ from the oracle waterfilling benchmark $\mathbf{D}_{\text{wf}}(\text{spec}(\Sigma_X), R)$?*

To make the comparison easier to interpret, we will phrase it in terms of rate overhead (of our codebook) compared to Σ_X -fine-tuned optimal codebook. Specifically, for a fixed spectrum $\lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ and a distortion level D^* , let $\mathbf{R}_{\text{wf}}(\lambda, D^*) = R_{\text{wf}}(\lambda, t)$ be the minimum oracle rate achieving $D_{\text{wf}}(\lambda, t) = D^*$ (see Eq. (WF)), and let $\mathbf{R}_{\text{rc}}(\lambda, D^*) = R_{\text{rc}}(\lambda, T)$ be the minimum random-coding rate achieving $D_{\text{rc}}(\lambda, T) = D^*$ (see Eq. (RDRC)). The difference $\mathbf{R}_{\text{rc}}(\lambda, D^*) - \mathbf{R}_{\text{wf}}(\lambda, D^*)$ is the rate overhead incurred by using a universal decoder g that is agnostic to the covariance matrix Σ_X with $\text{spec}(\Sigma_X) = \lambda$.

Our second main result in Theorem 5 shows that this overhead is uniformly bounded by 0.11 bit. Specifically, we derive precise expressions for $\mathbf{R}_{\text{rc}}(\lambda, D^*) - \mathbf{R}_{\text{wf}}(\lambda, D^*)$ that depend on $\lambda = \text{spec}(\Sigma_X)$ and D^* and prove that the maximum gaps occur at the spectra with at most 2 distinct eigenvalues and vanishing distortions (see Fig. 1 for the worst-case rate gap found at each $R = \mathbf{R}_{\text{rc}}(\lambda, D^*)$). See Sec. D for the full proof.

Theorem 5 (Main Result II: Worst-Case Rate Gap to Oracle Setting)

$$\sup_{D^* \in (0, 1)} \sup_{\Lambda} \{ \mathbf{R}_{\text{rc}}(\Lambda, D^*) - \mathbf{R}_{\text{wf}}(\Lambda, D^*) \} \leq 0.11,$$

where the supremum is over $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ with $\text{tr}(\Lambda) = n$.

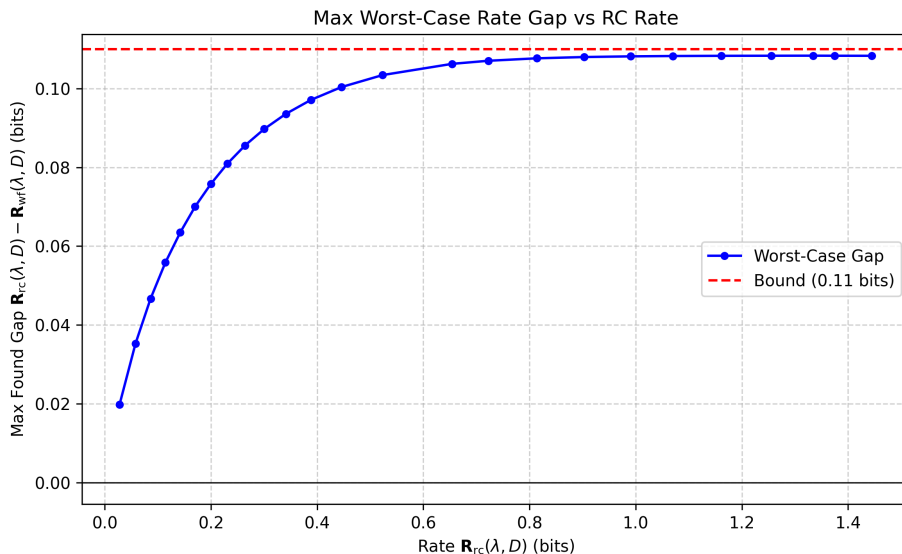


Figure 1: Maximum rate gap found at each rate R .

An intuitive way to see why the gap is bounded is to consider the most extreme of $\text{tr } \Sigma_X = n$ cases: the identity and the rank-1 case, for which $\mathbf{R}_{rc}(\lambda, D^*) = \mathbf{R}_{wf}(\lambda, D^*)$ by a simple computation. In fact, the gap vanishes for all matrices with semi-flat spectra $\lambda = (n/m, \dots, n/m, 0, \dots, 0)$, $m \leq n$, a phenomenon we further discuss in Sec. 3.1.1.

Discussion and Open Questions. Our results prove the *existence* of a universal codebook whose performance is uniformly within a constant rate gap of the Σ_X -aware (waterfilling) oracle benchmark, showing that universality is not, by itself, an information-theoretic bottleneck. Our random-coding-based proof is nonconstructive and does not yield an explicit codebook design or efficient encoder/decoder pair. Designing explicit and computationally efficient constructions is the most immediate open direction.

Often, low complexity quantizers are based on lattices. Indeed, for any fixed Σ_X a lattice randomly drawn from the natural Haar-Siegel measure will be a good quantizer with probability $1 - e^{-\Omega(n)}$. This follows from the results of [Ordentlich et al. \(2022\)](#) that show that the covering radius of a random lattice is typically near-optimal with respect to any (fixed) norm. However, *no lattice can be near-optimal simultaneously for all Σ_X* , and the reason is simple: for any fixed lattice $L \subset \mathbb{R}^n$ there is a rotation $U^\top L$ that aligns its directions with the natural basis. Consequently, for this rotation $U^\top L = \prod_{i=1}^n (\alpha_i \mathbb{Z})$ for some $\alpha_1, \dots, \alpha_n$. The integer lattice is a bad quantizer, and therefore any lattice quantizer must lose at least $\frac{1}{2} \log \frac{2\pi e}{12} \approx 0.254$ bits with respect to the waterfilling benchmark for some Σ_X (see [Ordentlich and Polyanskiy \(2026\)](#) for more details). It therefore follows that Σ_X -universal near optimal schemes cannot solely rely on lattice quantizers.

One example of a lattice-based algorithm that provides a practical solution to the problem studied here is the GPTQ algorithm [Frantar et al. \(2023\)](#) (with appropriate shaping/entropy coding [Ordentlich and Polyanskiy \(2026\)](#)). Its high-rate gap to the waterfilling benchmark, for particular Σ_X , is $\frac{1}{2} \log \frac{2\pi e}{12} + \frac{1}{2} \log(\text{AM} - \text{GM}(\Sigma_X))$ where $\text{AM} - \text{GM}(\Sigma_X)$ is the ratio between arithmetic-mean and geometric-mean of the squared diagonal elements in the Cholesky decomposition of Σ_X , which can be unbounded in general [Ordentlich and Polyanskiy \(2026\)](#). When W is a matrix consisting of $a \gg 1$ rows (rather than the vector case considered here), some of this gap can be reduced by send-

ing to the decoder $o(na)$ bits of information on Σ_X (which has negligible effect on the rate). One such example is the WaterSIC quantization scheme [Lifar et al. \(2026\)](#).

Related literature. The problem of vector quantization is classical [Gersho and Gray \(2012\)](#), and its asymptotic behavior (for iid sources and additive distortion) is given by the famous rate-distortion formulas, e.g. ([Polyanskiy and Wu, 2024](#), Part V). A recent wave of interest, however, focuses on computing quantized inner-product and matrix multiplication.

On the practical side, the pioneering work of [Dettmers et al. \(2022\)](#) demonstrated that substantial compression is possible via scaling plus uniform rounding (INT8 weight quantization), thus establishing the field of post-training quantization (PTQ). Notable PTQ works include SmoothQuant [Xiao et al. \(2024\)](#), which introduced calibration-based methods, i.e. those which depend on statistics of activations X via Σ_X . GPTQ [Frantar et al. \(2023\)](#) and LDLQ [Chee et al. \(2024\)](#), which are equivalent, simultaneously introduced an algorithm for Σ_X -dependent quantization. Going beyond simple integer-rounding, [Tseng et al. \(2024\)](#) consider lattice, [Savkin et al. \(2025\)](#) consider nested lattice and [Tseng et al. \(2025a\)](#) consider trellis quantization methods, respectively. Going beyond simple quadratic losses, are [Tseng et al. \(2025b\)](#) and [Badri and Shaji \(2023\)](#). We note that [Tseng et al. \(2024\)](#) also reintroduced random Hadamard transform (RHT) as a way of mitigating outliers, following earlier usage in quantization of gradients, and classically. See [Ashkboos et al. \(2024\)](#); [Liu et al. \(2025\)](#); [Chen et al. \(2025a\)](#) for other applications of RHT in PTQ.

On the theoretical side, the work [Ordentlich and Polyanskiy \(2025\)](#) established fundamental limits of quantized matrix multiplication by leveraging nested lattice quantization. The GPTQ/LDLQ algorithm was understood as Babai’s nearest-plane algorithm applied after a Cholesky factorization of Σ_X in [Chen et al. \(2025b\)](#); [Birnicks \(2025\)](#), and as a successive interference cancellation (SIC) algorithm in [Ordentlich and Polyanskiy \(2026\)](#), thus connecting weight-only quantization to lattice decoding and approximate closest vector problems [Conway and Sloane \(1982\)](#); [Babai \(1986\)](#). Authors of [Ordentlich and Polyanskiy \(2026\)](#) developed theoretical high-rate analysis of GPTQ, showed it can be arbitrarily far from waterfilling and proposed an improved algorithm which provably matches waterfilling to within 0.255 bit in the high-rate regime.

The question considered here (Σ_X -oblivious quantization) in the special case of diagonal Σ_X falls under the umbrella of the *compression with distortion as side-information* proposed in [Martinian et al. \(2008\)](#). This viewpoint connects modern LLM quantization to a long tradition in lattice decoding and approximate closest vector problems [Conway and Sloane \(1982\)](#); [Babai \(1986\)](#).

3. Technical Overview

As in any source-coding problem, once the codebook $\mathbf{C} = \{c_1, \dots, c_M\} \subset \mathbb{R}^n$ is fixed, the optimal encoder computes $i^* = i^*(W) = \operatorname{argmin}_{i \in [M]} d_{\Sigma_X}(W, c_i)$, and sends the index i^* to the decoder, which in turn outputs $\widehat{W} = c_{i^*}$. Thus, the distortion of the codebook \mathbf{C} is

$$nD(\mathbf{C}, \Sigma_X) = \mathbb{E} \left[\min_{i \in [M]} d_{\Sigma_X}(W, c_i) \right], \quad (4)$$

where the expectation is with respect to $W \sim \mathcal{N}(0, I_n)$. Even if we could design \mathbf{C} based on Σ_X , the distortion $D(\mathbf{C}, \Sigma_X)$ must satisfy (see [Sec. 2](#)) the waterfilling lower bound

$$D(\mathbf{C}, \Sigma_X) \geq \mathbf{D}_{\text{wf}}(R, \Sigma_X);$$

the lower bound is also asymptotically achievable in the limit of large n (see Proposition 2). The challenge is to find a single codebook \mathbf{C} with $M = 2^{nR}$ codewords in \mathbb{R}^n that attains small $D(\mathbf{C}, \Sigma_X) - \mathbf{D}_{\text{wf}}(R, \Sigma_X)$ simultaneously for all Σ_X .

As is standard, we prove the existence of such a codebook \mathbf{C} by drawing a random code with $M = 2^{nR}$ iid codewords from a distribution $P_{\widehat{W}}$. We show in Theorem 3 that for appropriate choice of $P_{\widehat{W}}$ we have that

$$\Pr_{\mathbf{C}} \left[\sup_{\Sigma_X} (D(\mathbf{C}, \Sigma_X) - \mathbf{D}_{\text{rc}}(\text{spec}(\Sigma_X), R - \varepsilon)) < \eta \right] \geq 1 - \exp(-\text{poly}(n)) \quad (5)$$

holds for any $\eta, \varepsilon > 0$ and n large enough, where $\text{spec}(\Sigma_X)$ is the vector of the eigenvalues of Σ_X and \mathbf{D}_{rc} is defined in Eq. (D_{rc}). Consequently there must exist a fixed rate- R codebook \mathbf{C} with

$$D(\mathbf{C}, \Sigma_X) \leq \mathbf{D}_{\text{rc}}(\text{spec}(\Sigma_X), R - \varepsilon) + \eta, \quad \forall \Sigma_X \in \mathbb{S}_+^n \text{ with } \text{tr}(\Sigma_X) = n. \quad (6)$$

3.1. Codebook Distribution

How should we choose $P_{\widehat{W}}$? For a given Σ_X with spectral decomposition $\Sigma_X = U\Lambda U^\top$, the optimal $P_{\widehat{W}}$ follows from the waterfilling solution. Specifically, for water-level $1/t$ chosen so that R_{wf} defined in (WF) equals R , the optimal distribution is

$$P_{\widehat{W}}^*(\Sigma_X, R) = P_{\widehat{W}}^*(\Sigma_X, t) = \mathcal{N}(0, U \Gamma(\Lambda, t) U^\top),$$

$$\text{where } \Gamma(\Lambda, t) = \text{diag} \left(\max \left\{ 1 - \frac{t}{\lambda_1}, 0 \right\}, \dots, \max \left\{ 1 - \frac{t}{\lambda_n}, 0 \right\} \right). \quad (7)$$

We need to choose a single $P_{\widehat{W}}$ that “works well” for all Σ_X . Since there is no preference to any $U \in \mathcal{O}_n$, it makes sense to take an isotropic Gaussian $P_{\widehat{W}}$. Observing that for any fixed Σ_X and $t > 0$ the covariance matrix for $P_{\widehat{W}}^*(\Sigma_X, t)$ satisfies $U \Gamma(\Lambda, t) U^\top \preceq I_n$, we will take

$$P_{\widehat{W}}(\tau) = \mathcal{N}(0, \tau^2 I_n), \quad (8)$$

for some $0 < \tau < 1$.

An appealing feature of the isotropic Gaussian distribution is that drawing M iid vectors from $\mathcal{N}(0, \tau^2 I_n)$ is equivalent to first drawing them iid from $\mathcal{N}(0, I_n)$ and then scaling all of them by τ . The consequence of this simple fact is that while we cannot perfectly match our codebook distribution to $P_{\widehat{W}}^*(\Sigma_X, R)$, the flexibility in the choice of $\tau = \tau(\Sigma_X, R)$ allows for a better match. Consequently, we draw the M codewords of \mathbf{C} from the $\mathcal{N}(0, I_n)$ distribution. The encoder, that knows Σ_X , computes $\tau(\Sigma_X, R)$ that provides the smallest expected distortion, and sends a description of τ to the decoder.² With this procedure, the effective codebook $\tilde{\mathbf{C}} = \tau \mathbf{C}$ is drawn from $P_{\widehat{W}} = \mathcal{N}(0, \tau^2(\Sigma_X, R)I_n)$. The encoder then finds $i^* = \text{argmin}_{i \in [M]} d_{\Sigma_X}(W, \tilde{c}_i) = \text{argmin}_{i \in [M]} d_{\Sigma_X}(W, \tau c_i)$ and sends i^* as well as a high-resolution description of τ in bits to the decoder.

2. As we will see below, some further gain can be attained by allowing τ to also depend on the source realization $w \in \mathbb{R}^n$.

3.1.1. GEOMETRIC INTUITION

To get some intuition to why the “universal” codebook distribution $P_{\widehat{W}} = \mathcal{N}(0, \tau^2(\Sigma_X, R)I_n)$ works well simultaneously for all Σ_X , let us restrict attention to the family of covariance matrices with semi-flat spectrum. In particular, for $m \leq n$ let

$$\mathcal{S}_m^n = \left\{ \Sigma_X = U\Lambda U^\top : U \in \mathcal{O}_n, \lambda_1 = \dots = \lambda_m = \frac{n}{m}, \lambda_{m+1} = \dots = \lambda_n = 0 \right\}, \quad (9)$$

be the collection of PSD matrices with $m \leq n$ equal and non-zero singular values and $m - n$ zero singular values. From (7) we see that the Σ_X -matched optimal codebook distribution is of the form

$$P_{\widehat{W}}^*(\Sigma_X, R) = \mathcal{N}\left(0, U \cdot \text{diag}(\tau^2, \dots, \tau^2, 0, \dots, 0) U^\top\right) \quad (10)$$

for some $\tau > 0$. Thus the optimal procedure for random coding is to draw iid Gaussian codewords within the subspace $U_{[m]}$ spanned by the first m singular vectors. Our universal distribution, on the other hand, draws isotropic iid Gaussian codewords, and is hence very far from the optimal distribution. However, since the encoder searches for the nearest codeword under the d_{Σ_X} metric, it effectively projects both W and the codebook \mathbf{C} to $U_{[m]}$ and finds the nearest codeword in ℓ_2 -metric within this subspace. It therefore follows that what dictates performance of a random code \mathbf{C} under d_{Σ_X} metric (for Σ_X with semi-flat spectrum) is the distribution of $U_{[m]}^\top c$, where $c \sim P_{\widehat{W}}$. Thus, our universal $P_{\widehat{W}} = \mathcal{N}(0, \tau^2 I_n)$ is simultaneously optimal for all semi-flat Σ_X (with all possible $m \in [n]$), provided that we judiciously choose $\tau = \tau(\Sigma_X, R)$. Inspection of our Σ_X -universal rate-distortion tradeoff (RDRC) shows that indeed

$$\mathbf{D}_{\text{rc}}(\text{spec}(\Sigma_X), R) = \mathbf{D}_{\text{wf}}(\text{spec}(\Sigma_X), R),$$

for all semi-flat Σ_X . Whenever the spectrum of Σ_X is not semi-flat the distribution of $U^\top c$ for $c \sim P_{\widehat{W}}$ does not match that of $U^\top c$ under the optimal $P_{\widehat{W}}^*(\Sigma_X, R)$, and consequently in these cases our Σ_X -universal rate-distortion tradeoff (RDRC) is worse than the waterfilling rate-distortion tradeoff. Nevertheless, it turns out that the loss for this mismatch is at most 0.11 bits, as shown in Theorem 5.

3.2. Sketch of Proof

After we obtained intuition for the choice of using a random iid isotropic Gaussian codebook, with scale τ determined by the encoder, we move on to giving an overview of the proof of Theorem 3. The detailed rigorous proof is given in Appendix C.

Let us first fix $\Sigma_X = U\Lambda U^\top$, and analyze the performance of a codebook $\mathbf{C} = \{c_1, \dots, c_M\}$, $M = 2^{nR}$, with $c_i \stackrel{iid}{\sim} \mathcal{N}(0, I_n)$. We will show that \mathbf{C} is “good” for Σ_X with probability $1 - \exp(-e^{\Omega(n)})$, and from this we will deduce that a random \mathbf{C} is “good” for all Σ_X by a covering argument.

The codebook \mathbf{C} can describe a fixed $w \in \mathbb{R}^n$ with distortion $\leq D$ if at least one of its τ -scaled codewords is inside the region $w + \sqrt{D}\mathcal{B}_{\Sigma_X}$, where

$$\mathcal{B}_{\Sigma_X} = \left\{ e \in \mathbb{R}^n : e^\top \Sigma_X e \leq n \right\}. \quad (11)$$

Thus, the key to analyzing the tradeoff between rate and distortion is understanding how the success probability of a single codeword behaves as a function of D . Since the codewords are $\mathcal{N}(0, I_n)$, their distribution is invariant to rotation, and therefore the success probability of a single codeword is

$$\begin{aligned} p_{\text{success}}(w, \tau, \Sigma_X, D) &= \phi_{\tau^2} \left(w + \sqrt{D} \mathcal{B}_{\Sigma_X} \right) \\ &= \phi_{\tau^2} \left(U^\top w + \sqrt{D} \mathcal{B}_\Lambda \right) = p_{\text{success}}(U^\top w, \tau, \Lambda, D), \end{aligned} \quad (12)$$

where ϕ_{τ^2} denotes the probability distribution for $\mathcal{N}(0, \tau^2 I_n)$. When $p_{\text{success}} < 2^{-n(R+\varepsilon)}$, it is very unlikely to find a codeword in $w + \sqrt{D} \mathcal{B}_{\Sigma_X}$, and on the other hand, when $p_{\text{success}} > 2^{-n(R-\varepsilon)}$ we are very likely to find a codeword in $w + \sqrt{D} \mathcal{B}_{\Sigma_X}$. Thus, what we are looking for is the critical D for which $-\frac{1}{n} \log p_{\text{success}}(U^\top w, \tau, \Lambda, D) \approx R$. Since the codebook's scale τ needs to be sent from the encoder to the decoder anyway, we may let it depend not only on Λ but also on $U^\top w$. Therefore, given $U^\top w, \Lambda$ and R we choose $\tau = \tau(U^\top w, \Lambda, R)$ for which the critical D is small. A tedious but straightforward calculation shows that the optimal choice is

$$\tau = \tau(U^\top w, \Lambda, R) = \left(T \sum_j \frac{(U^\top w)_j^2 \lambda_j^2}{(1 + \lambda_j T)^2} \right)^{1/2} \left(\sum_j \frac{\lambda_j}{1 + \lambda_j T} \right)^{-1/2}, \quad (13)$$

where $T = T(\Lambda, R)$ is such that $R_{\text{rc}}(\Lambda, T) = R$, and $R_{\text{rc}}(\Lambda, T)$ is defined in (RDRC). Let

$$D_{\text{rc}}(U^\top w) = D_{\text{rc}}(\Lambda, R, U^\top w) = \frac{1}{n} \sum_{i=1}^n \frac{(U^\top w)_i^2 \lambda_i}{1 + \lambda_i T}, \quad (14)$$

where here as well $T = T(\Lambda, R)$. In Lemma 10 we prove that

$$-\log p_{\text{success}}(U^\top w, \tau, \Lambda, D_{\text{rc}}(U^\top w)) = -\log \phi_{\tau^2} \left(U^\top w + \sqrt{D_{\text{rc}}(U^\top w)} \mathcal{B}_\Lambda \right) \approx nR. \quad (15)$$

Thus, $D_{\text{rc}}(U^\top w)$ is the critical distortion for fixed $w \in \mathbb{R}^n$ and Σ_X . Note that for $W \sim \mathcal{N}(0, I_n)$ we have

$$\mathbb{E}_W \left[D_{\text{rc}}(U^\top W) \right] = \frac{1}{n} \sum_{i=1}^n \frac{\lambda_i}{1 + \lambda_i T} = D_{\text{rc}}(\Lambda, R), \quad (16)$$

where $D_{\text{rc}}(\Lambda, R)$ is the distortion in (D_{rc}).

The proof of Lemma 10 uses standard large deviation techniques, but note that what we really need for the analysis is a quantitative lower bound on p_{success} , and this requires some more work beyond Chernoff bound. The choice of τ from (13) is the one that minimizes the large deviations exponent, whereas the parameter T in our rate-distortion tradeoff (RDRC) is just a rescaling of the parameter t in Chernoff's bound $\Pr(X > D) \leq e^{-tD} \mathbb{E}[e^{tX}]$.

By (15), for $\eta > 0$ we have that $p_{\text{success}}(U^\top w, \tau, \Lambda, D_{\text{rc}}(U^\top w) + \eta/2) \geq 2^{-n(R-\varepsilon_\eta)}$ for some $\varepsilon_\eta > 0$ (and n large enough). Thus, $\forall w \in \mathbb{R}^n$

$$\begin{aligned} P_{\text{failure}}(w) &= \Pr_{\mathbf{C}} \left(\frac{1}{n} \min_{i \in [M]} d_{\Sigma_X}(w, \tau c_i) > D_{\text{rc}}(U^\top w) + \eta/2 \right) \\ &= 1 - \left(1 - p_{\text{success}} \left(U^\top w, \tau, \Lambda, D_{\text{rc}}(U^\top w) + \eta/2 \right) \right)^M \leq \exp(-2^{n\varepsilon}), \end{aligned} \quad (17)$$

for some $\varepsilon > 0$. For a *fixed* code \mathbf{C} define

$$\mathcal{E}_{\text{failure}}(\mathbf{C}) = \left\{ w \in \mathbb{R}^n : \frac{1}{n} \min_{i \in [M]} d_{\Sigma_X}(w, \tau c_i) > D_{\text{rc}}(U^\top w) + \eta/2 \right\}. \quad (18)$$

Assuming we can always encode w to 0 (e.g., by setting $\tau = 0$ if needed), for any $W \sim P_W$ we have

$$\begin{aligned} \frac{1}{n} D(\mathbf{C}, \Lambda) &= \mathbb{E} \left[\min_{i \in [M]} d_{\Sigma_X}(W, \tau c_i) \right] \leq \mathbb{E} \left[D_{\text{rc}}(U^\top W) + \frac{\eta}{2} + \mathbb{1}\{W \in \mathcal{E}_{\text{failure}}(\mathbf{C})\} d_{\Sigma_X}(W, 0) \right] \\ &\leq \mathbb{E} \left[D_{\text{rc}}(U^\top W) \right] + \frac{\eta}{2} + \sqrt{P_W[\mathcal{E}_{\text{failure}}(\mathbf{C})]} \cdot \sqrt{\mathbb{E}[d_{\Sigma_X}^2(W, 0)]}, \end{aligned} \quad (19)$$

where the last inequality is Cauchy-Schwarz. If $\mathbb{E}_W[d_{\Sigma_X}^2(W, 0)] \leq \text{poly}(n)$ under $W \sim P_W$,³ it follows that whenever $P_W[\mathcal{E}_{\text{failure}}(\mathbf{C})]$ is sufficiently small, say $\leq e^{-n}$, the codebook \mathbf{C} attains the desired

$$\frac{1}{n} D(\mathbf{C}, \Lambda) \leq \mathbb{E} \left[D_{\text{rc}}(U^\top W) \right] + \eta.$$

$P_W[\mathcal{E}_{\text{failure}}(\mathbf{C})]$ can indeed be bounded: using (17) and Markov's inequality we obtain that this holds for the vast majority of codebooks:

$$\begin{aligned} \Pr_{\mathbf{C}} \left(P_W[\mathcal{E}_{\text{failure}}(\mathbf{C})] > e^{-n} \right) &\leq e^n \mathbb{E}_{\mathbf{C}} [P_W[\mathcal{E}_{\text{failure}}(\mathbf{C})]] = e^n \mathbb{E}_{\mathbf{C}, W} [P_{\text{failure}}(W)] \\ &\leq \exp(-2^{n\varepsilon} + n). \end{aligned} \quad (20)$$

From this we conclude that for fixed P_W and any fixed Σ_X

$$\Pr_{\mathbf{C}} \left(\left[\frac{1}{n} D(\mathbf{C}, \Lambda) - \mathbb{E}_W \left[D_{\text{rc}}(U^\top W) \right] > \eta \right] \right) \leq \exp(-2^{n\varepsilon'}). \quad (21)$$

Specializing this to $P_W = \mathcal{N}(0, I_n)$ we obtain

$$\Pr_{\mathbf{C}} \left(\left[\frac{1}{n} D(\mathbf{C}, \Lambda) - D_{\text{rc}}(\Lambda, R) \right] > \eta \right) \leq \exp(-2^{n\varepsilon'}). \quad (22)$$

We have therefore obtained that for any fixed Σ_X the probability of drawing a “bad” codebook with $\left[\frac{1}{n} D(\mathbf{C}, \Lambda) - D_{\text{rc}}(\Lambda, R) \right] > \eta$ is double-exponentially small. From here, it is clear how to prove that a randomly drawn \mathbf{C} will have

$$\Pr_{\mathbf{C}} \left(\sup_{\Sigma_X} \left[\frac{1}{n} D(\mathbf{C}, \Lambda_X) - D_{\text{rc}}(\Lambda, R) \right] > \eta \right) < \exp(-2^{n\varepsilon''}). \quad (23)$$

All we need is to find a dense cover of PSD matrices in $R^{n \times n}$ with trace n whose size is $\exp(-2^{o(n)})$. In particular, we need to find a collection of $N = \exp(-2^{o(n)})$ PSD matrices with the property that for any valid PSD matrix Σ_X there exists a matrix $\Sigma_i, i \in [N]$ such that:

$$D_{\text{rc}}(\text{spec}(\Sigma_X), R) \approx D_{\text{rc}}(\text{spec}(\Sigma_i), R) \quad \text{and} \quad \frac{1}{n} D(\mathbf{C}, \Sigma_X) \approx \frac{1}{n} D(\mathbf{C}, \Sigma_i). \quad (24)$$

3. e.g., $\mathbb{E}_W[d_{\Sigma_X}^2(W, 0)] \leq O(n^2)$ for $W \sim \mathcal{N}(0, I_n)$

Since N is allowed to be so large, these two constraints can be met to arbitrary resolution (though the proof for this requires a lot of bookkeeping and is somewhat technical, see Sec. C.2).

We end this overview with listing some of the technical issues that our sketch of proof above ignored, and briefly mention how addressing them, as we of course do in the actual proofs, affects the results.

τ quantization. While our sketch assumed that τ can be conveyed to the decoder in perfect resolution, in reality, some of our nR bits budget is allocated to the description of τ . In order to compress τ that depends on $U^\top W$, we require a high probability bound on $\|U^\top W\|_\infty$. This norm constraint is also needed for the proof of our large deviations result (Lemma 10).

Norm bound for codewords. The perturbation argument (24) for the dense cover must account for quantization and change in τ . Since these errors are multiplied by codewords from \mathbf{C} , we require a uniform norm bound $\|c_i\|_2 < n^B$ for some $B > 1$ (say $B = 10$) to make their contribution to end-to-end distortion negligible. While the probability that this occurs is overwhelmingly large, it is only $1 - \exp(-\text{poly}(n))$ rather than double exponential. For this reason the probability of drawing a Σ_X -universal codebook is only $1 - \exp(-\text{poly}(n))$ (see Eq. (5)) rather than the $1 - \exp(-2^{\Omega(n)})$ that our sketch of proof gives. Because event $\{\|c_i\|_2 < n^B \ \forall i \in [M]\}$ is independent of Σ_X , it contributes only a single term to the union bound. Consequently, our dense grid is still allowed to be of size $\exp(-2^{o(n)})$.

Rate-penalty. In the overview above, we assumed that if $-\ln p_{\text{success}}(D) = R$, then for any $\eta > 0$ we have $-\ln p_{\text{success}}(D + \eta) = R - \varepsilon_\eta$ for some $\varepsilon_\eta > 0$. This is indeed the case whenever $D'(R)$ is finite. However, our results are for the supremum over all Σ_X with trace n , and this cannot be guaranteed for all such matrices at all rates. Consequently, in Theorem 3 there is both a rate-penalty $\varepsilon > 0$ and a distortion penalty $\eta > 0$ (which both can be made arbitrarily small for n large enough), whereas in the sketch above we only had a distortion penalty.

Disclosure of LLM Assistance

The authors used ChatGPT to assist with editing, code generation for Fig. 1, and technical steps (e.g., perturbation analysis and derivative computations) in the proofs of main theorems. The final optimization step after the spectrum reduction in Theorem 5 was also suggested by ChatGPT. All model-generated code, computations, and proof steps were independently verified and adjusted by the authors. The authors take full responsibility for the correctness of all analytical and numerical results in the paper.

References

- Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L. Croci, Bo Li, Pashmina Cameron, Martin Jaggi, Dan Alistarh, Torsten Hoeffler, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms, 2024. URL <https://arxiv.org/abs/2404.00456>.
- László Babai. On lovász’lattice reduction and the nearest lattice point problem. *Combinatorica*, 6(1):1–13, 1986.
- Hicham Badri and Appu Shaji. Half-quadratic quantization of large machine learning models, November 2023. URL https://mobiusml.github.io/hqq_blog/.

- Johann Birnick. The lattice geometry of neural network quantization—a short equivalence proof of gptq and babai’s algorithm. *arXiv preprint arXiv:2508.01077*, 2025.
- Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher De Sa. Quip: 2-bit quantization of large language models with guarantees, 2024. URL <https://arxiv.org/abs/2307.13304>.
- Jiale Chen, Vage Egiazarian, Torsten Hoefer, and Dan Alistarh. Wush: Near-optimal adaptive transforms for llm quantization. *arXiv preprint arXiv:2512.00956*, 2025a.
- Jiale Chen, Yalda Shabanzadeh, Elvir Crnčević, Torsten Hoefer, and Dan Alistarh. The geometry of llm quantization: Gptq as babai’s nearest plane algorithm. *arXiv preprint arXiv:2507.18553*, 2025b.
- John Conway and Neil Sloane. Fast quantizing and decoding and algorithms for lattice quantizers and codes. *IEEE Transactions on Information Theory*, 28(2):227–232, 1982.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35: 30318–30332, 2022.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefer, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers, 2023. URL <https://arxiv.org/abs/2210.17323>.
- Allen Gersho and Robert M Gray. *Vector quantization and signal compression*, volume 159. Springer Science & Business Media, 2012.
- Egor Lifar, Semyon Savkin, Or Ordentlich, and Yury Polyanskiy. Watersic: information-theoretically (near) optimal linear layer quantization. *arXiv preprint arXiv:2603.04956*, 2026.
- Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. Spinquant: Llm quantization with learned rotations, 2025. URL <https://arxiv.org/abs/2405.16406>.
- Emin Martinian, Gregory W. Wornell, and Ram Zamir. Source coding with distortion side information. *IEEE Transactions on Information Theory*, 54(10):4638–4665, 2008. doi: 10.1109/TIT.2008.928983.
- NVIDIA et al. Pretraining large language models with NVFP4. *arXiv preprint arXiv:2509.25149*, 2025.
- Open Compute Project. OCP microscaling formats (MX) specification. Technical report, Open Compute Project, 2023. URL <https://www.opencompute.org/documents/ocp-microscaling-formats-mx-v1-0-spec-final-pdf>.
- Or Ordentlich and Yury Polyanskiy. Optimal quantization for matrix multiplication. *IEEE Transactions on Information Theory*, 2025.
- Or Ordentlich and Yury Polyanskiy. High-rate quantized matrix multiplication II. *arXiv preprint arXiv:2605.13768*, 2026.

- Or Ordentlich, Oded Regev, and Barak Weiss. New bounds on the density of lattice coverings. *Journal of the American Mathematical Society*, 35(1):295–308, 2022.
- Jan Ostergaard and Ram Zamir. Incremental refinement using a gaussian test channel. In *2011 IEEE International Symposium on Information Theory Proceedings*, pages 2233–2237, 2011. doi: 10.1109/ISIT.2011.6033957.
- Yury Polyanskiy and Yihong Wu. *Information theory: From coding to learning*. Cambridge university press, 2024.
- Semyon Savkin, Eitan Porat, Or Ordentlich, and Yury Polyanskiy. Nestquant: Nested lattice quantization for matrix products and llms, 2025. URL <https://arxiv.org/abs/2502.09720>.
- Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even better llm quantization with hadamard incoherence and lattice codebooks, 2024. URL <https://arxiv.org/abs/2402.04396>.
- Albert Tseng, Qingyao Sun, David Hou, and Christopher De Sa. Qtip: Quantization with trellises and incoherence processing, 2025a. URL <https://arxiv.org/abs/2406.11235>.
- Albert Tseng, Zhaofeng Sun, and Christopher De Sa. Model-preserving adaptive rounding, 2025b. URL <https://arxiv.org/abs/2505.22988>.
- Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models, 2024. URL <https://arxiv.org/abs/2211.10438>.
- Ram Zamir. The rate loss in the wyner-ziv problem. *IEEE Transactions on Information Theory*, 42(6):2073–2084, 2002.
- Ram Zamir and Toby Berger. Multiterminal source coding with high resolution. *IEEE Transactions on Information Theory*, 45(1):106–117, 2002.

Appendix A. Preliminaries and Notation

Notation. We write $\Sigma \succeq 0$ to denote that the matrix $\Sigma \in \mathbb{R}^{n \times n}$ is positive semidefinite; we define the set of positive semidefinite matrices as $\mathbb{S}_+^n = \{\Sigma \in \mathbb{R}^{n \times n} : \Sigma \succeq 0\}$.

For $\Sigma \in \mathbb{S}_+^n$ we denote the spectral decomposition as

$$\text{EVD}(\Sigma) = U\Lambda U^\top,$$

where $U \in \mathcal{O}_n$ is orthogonal and Λ is diagonal.

$\text{diag}(v)$ for $v \in \mathbb{R}^n$ denotes an $n \times n$ diagonal matrix with $\text{diag}(v)_{ii} = v_i$. $\text{spec}(A)$ for $A \in \mathbb{R}^{n \times n}$ denotes the vector of the eigenvalues of A .

For $A \in \mathbb{R}^{n \times a}$ and $\Sigma \in \mathbb{S}_+^n$, denote

$$\|A\|_\Sigma \triangleq \sqrt{\text{tr}(A^\top \Sigma A)}.$$

Proposition 6 (Hanson-Wright Concentration Inequality) For $X \sim \mathcal{N}(0, I_n)$ and $A \in \mathbb{R}^{n \times n}$, for every $t \geq 0$,

$$\Pr \left[|X^\top A X - \text{tr} A| > t \right] \leq 2 \exp \left(-c \min \left(\frac{t^2}{K^4 \|A\|_F^2}, \frac{t}{K^2 \|A\|_{op}} \right) \right),$$

for universal constants c, K .

Appendix B. Random Coding: Worst-Case W

In this section, for fixed $\Sigma_X \in \mathbb{S}_+^n$, we characterize the rate-distortion of quantizing a fixed vector $W \in \mathbb{R}^n$ under a distortion metric d_{Σ_X} using a random coding scheme. In particular, for a given $\Sigma_X = U^\top \Lambda U$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$, fixed vector $W \in \mathbb{R}^n$, and constant $R^* > 0$, we define a generalized distortion function $\mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*)$ and demonstrate that it is achievable in Theorem 7.

General Rate-Distortion Function. Let $V, \lambda = (\lambda_1, \dots, \lambda_n)^\top \in \mathbb{R}^n$ be such that $\lambda_i \geq 0$ for all $i \in [n]$ and $\sum_i \lambda_i = n$. We define *general random-coding rate-distortion function* in dimension n to be the following parametric curve for $T > 0$:

$$D_{\text{rc}}^\lambda(V, T) = \frac{1}{n} \sum_{i=1}^n \frac{V_i^2 \lambda_i}{1 + \lambda_i T} \quad \text{and} \quad R_{\text{rc}}^\lambda(T) = \frac{1}{2n} \sum_{i=1}^n \log(1 + \lambda_i T). \quad (\text{RDRC})$$

Throughout, \log denotes \log_2 and \ln denotes the natural logarithm. Denote $T_{\text{rc}}^\lambda(R)$ be a unique value T , s.t. $R_{\text{rc}}^\lambda(T) = R$ (note that $T_{\text{rc}}^\lambda(R)$ is independent of V). Let

$$\mathbf{D}_{\text{rc}}^\lambda(V, R) \triangleq D_{\text{rc}}^\lambda(V, T_{\text{rc}}^\lambda(R)). \quad (\mathbf{D}_{\text{rc}})$$

We consider a task of quantizing a given vector $W \in \mathbb{R}^n$ under a distortion function (d_{Σ_X}).

Condition 1 (Admissible W, Σ_X) We consider $W \in \mathbb{R}^n, \Sigma_X \in \mathbb{S}_+^n$ satisfying

1. $\text{EVD}(\Sigma_X) = U\Lambda U^\top$ for $U \in \mathcal{O}_n$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ with $\text{tr}(\Lambda) = n$;
2. $\|U^\top W\|_\infty \leq n^\alpha$ for a known constant α .

In Theorem 7 we show that there exists an encoder-decoder pair that, for any constant target rate $R^* > 0$ and fixed admissible W, Σ_X (Cond. 1), achieves rate R^* and distortion $\mathbf{D}_{\text{rc}}^{U^\top W, \lambda}(R^*)$ asymptotically, with high probability over the randomness S shared between encoder and decoder (from which a codebook is generated).⁴

Theorem 7 (Achievability of Random-Coding Rate-Distortion: Nonasymptotic Guarantee) *Fix any constants $R^*, \varepsilon_o, \eta_o > 0$ and $\alpha > 0$. There exists an encoder $f : \mathbb{R}^n \times \mathbb{S}_+^n \times [0, 1] \rightarrow [2^{nR}]$, a decoder $g : [2^{nR}] \times [0, 1] \rightarrow \mathbb{R}^n$ with $R \leq R^* + \varepsilon_o$ and a (shared) random variable $S \in [0, 1]$ with the following property. For any fixed W, Σ_X (satisfying Cond. 1) we set $\widehat{W} = g(f(W, \Sigma_X, S), S)$. Then for any sufficiently large $n \geq n_0 = n_0(\varepsilon_o, \eta_o, R^*, \alpha)$ and any $\beta > \alpha - 1/4$ in case $\alpha \geq 1/4$ and $\beta = 0$ otherwise, we have*

$$\Pr_S \left[\frac{1}{n} d_{\Sigma_X}(W, \widehat{W}) \leq \mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*) + n^{2\beta} \eta_o \right] \geq 1 - \exp \left(-2^{n\varepsilon_o(1-cn^{4(\alpha-\beta)-1})} \right),$$

where $d_{\Sigma_X}(W, \widehat{W})$ is the distortion function in Eq. (d $_{\Sigma_X}$) and $c = c(\varepsilon_o, \eta_o, R^*, \alpha)$.

Remark 8 *In case of $\alpha < 1/4$, the distortion bound above simplifies to $\mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*) + \eta$ (since $\beta = 0$).*

B.1. Proof of Theorem 7

Quantization Scheme. From the shared randomness S , f, g generate a Gaussian codebook $\mathbf{C} = \{c_1, \dots, c_{M=2^{n(R^*+\varepsilon)}}\}$ for $\varepsilon = \varepsilon(\varepsilon_o)$ to be chosen later. Denote $\widetilde{W} = U^\top W$. We define:

- **Encoder f :** Let $T := T_{\text{rc}}^\lambda(R^*)$ and define the scaling parameter $\tau = \tau(\widetilde{W}, \lambda)$ as

$$\tau = \left(T \sum_j \frac{\widetilde{W}_j^2 \lambda_j^2}{(1 + \lambda_j T)^2} \right)^{1/2} \left(\sum_j \frac{\lambda_j}{1 + \lambda_j T} \right)^{-1/2}. \quad (\tau \text{ def.})$$

We set f to be a tuple

$$f(W, \Sigma_X, S) = \left(\underset{i \in [M]}{\operatorname{argmin}} d_{\Sigma_X}(W, \tau c_i), q(\tau) \right),$$

where $q(\tau) = \delta \lceil \|\widetilde{W}\|_\infty \tau / (\delta \|\widetilde{W}\|_\infty) \rceil$ is a rounding quantization scheme of precision δ and recall the distortion function

$$d_{\Sigma_X}(W, C) := (W - C)^\top \Sigma_X (W - C).$$

- **Decoder g :**

$$g(i, q(\tau), S) = q(\tau) \cdot c_i.$$

4. In the regime $\alpha < 1/4$, see the statement of Thm. 7.

Rate-Distortion Bound. In the quantization scheme above,

$$R = \underbrace{R^* + \varepsilon}_{\text{Gauss. codebook}} + \underbrace{\frac{1}{n} \log(1/\delta)}_{\tau \text{ quant.}} .$$

The rest of the proof is to obtain a high probability bound on the resulting distortion that, given $\widehat{W} = g(f(W, \Sigma_X, S), S) = q(\tau) \cdot c_i$, can be expressed as:

$$d_{\Sigma_X}(W, \widehat{W}) = d_{\Sigma_X}(W, q(\tau)c_i) = (W - q(\tau)c_i)^\top \Sigma_X (W - q(\tau)c_i) .$$

In what follows, denote $D^* := \mathbf{D}_{\text{rc}}^\lambda(\widehat{W}, R^*)$.

Before giving the proof we state two helpful claims that bound the effect of quantizing τ in the scheme above. Proofs of Claim 1 and 2 are found in this subsection below.

Claim 1 (Bound on τ .) *The value τ in Eq. ([\tau def.](#)) satisfies*

$$0 \leq \tau \leq \|\widehat{W}\|_\infty \leq n^\alpha .$$

Claim 2 (Bound on distortion from τ quantization) *Given $|q(\tau) - \tau| = \delta_\tau \leq \delta n^\alpha$, with probability at least $1 - 2^{n(R^* + \varepsilon)} \exp(-Ct)$ for a universal constant C and any $t > 1/n$,*

$$d_{\Sigma_X}(W, q(\tau)c_i) \leq \left(\sqrt{d_{\Sigma_X}(W, \tau c_i)} + \delta n^\alpha \sqrt{n(1+t)} \right)^2 .$$

The main part of the argument is essentially contained in the following Lemma 9. The proof, which also explains the expressions for D_{rc} and R_{rc} is proven in a separate section [B.2](#) due to its importance.

Lemma 9 (Gaussian Book Success) *Let $\varepsilon, \eta > 0$ be constants and $\mathbf{C} = \{c_1, \dots, c_{M=2^{n(R^* + \varepsilon)}}\}$, $c_i \sim_{i.i.d.} \mathcal{N}(0, I_n)$ be a randomly generated Gaussian codebook.*

For admissible W, Σ_X (Cond. [1](#)) with $\alpha < 1/4$,

$$\Pr_{\mathbf{C}} \left[\frac{1}{n} \min_{i \in [M]} d_{\Sigma_X}(W, \tau(W, \Sigma_X) \cdot c_i) \leq \mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*) + \eta \right] \geq 1 - \exp\left(-2^{n\varepsilon(1-cn^{4\alpha-1})}\right) ,$$

where $c = c(\varepsilon, \eta, R^, \alpha)$ is an explicit constant function, $\mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*)$ is defined in Eq. ([D_{rc}](#)), and $\tau(W, \Sigma_X)$ is defined in Eq. ([\tau def.](#)).*

Proof (of Theorem [7](#)) We apply Lemma 9 to $W' = n^{-\beta}W$ for $\eta = \eta(\eta_0)$ to be chosen later and any constant $\beta > \alpha - 1/4$ in case $\alpha \geq 1/4$ and $\beta = 0$ otherwise. We have $\|U^\top W'\|_\infty = n^{-\beta}\|U^\top W\|_\infty \leq n^{\alpha-\beta} = o(n^{1/4})$, and therefore,

$$\begin{aligned} \Pr_{\mathbf{C}} \left[\min_{i \in [M]} d_{\Sigma_X}(W, \tau c_i) \leq n \left(\underbrace{\mathbf{D}_{\text{rc}}^\lambda(\widehat{W}, R^*)}_{D^*} + n^{2\beta} \eta \right) \right] &= \Pr_{\mathbf{C}} \left[\min_{i \in [M]} d_{\Sigma_X}(W', \tau n^{-\beta} c_i) \leq n \left(\mathbf{D}_{\text{rc}}^\lambda(U^\top W', R^*) + \eta \right) \right] \\ &\geq 1 - \exp\left(-2^{n\varepsilon(1-cn^{4(\alpha-\beta)-1})}\right) . \end{aligned}$$

From Claim 1, the simple rounding quantizer $q(\tau) = \delta \|\widetilde{W}\|_\infty \lceil \tau / (\delta \|\widetilde{W}\|_\infty) \rceil$ achieves $|q(\tau) - \tau| \leq \delta n^\alpha$, and therefore, by Claim 2 and a union bound, with probability at least $1 - \exp\left(-2^{n\varepsilon(1-cn^4(\alpha-\beta)^{-1})}\right) - 2^{n(R^*+\varepsilon)} \exp(-Ct)$,

$$d_{\Sigma_X}(W, q(\tau)c_i) \leq \left(\sqrt{n(D^* + n^{2\beta}\eta)} + \delta n^\alpha \sqrt{n(1+t)} \right)^2. \quad (25)$$

It remains to simplify the expression in Eq. (25). Let $A > 0$ be any constant and set

$$\delta =: \min \left\{ \frac{1}{4\sqrt{1+t}} \left(\frac{1}{2\ln 2 \cdot R^*} + \eta \right)^{-1/2} n^{-A-2\alpha-1}, \frac{1}{\sqrt{2(1+t)}} n^{-(A+2\alpha+1)/2} \right\}. \quad (26)$$

Notice that since $\forall i, \frac{\lambda_i}{1+2T\lambda_i} \leq \frac{1}{2T}$, we have $D^* \leq \frac{n^{2\alpha}}{2T} \leq \frac{n^{2\alpha}}{2\ln 2 \cdot R^*}$, where the last inequality is derived by $\ln 2 \cdot R^* = \frac{1}{2n} \sum_i \ln(1 + 2\lambda_i T) \leq T$. Plugging in the δ value in Eq. (26) into Eq. (25), we obtain

$$\begin{aligned} d_{\Sigma_X}(W, q(\tau)c_i) &\leq n \left(D^* + n^{2\beta}\eta + 2\delta \sqrt{D^* + n^{2\beta}\eta} \cdot n^\alpha \sqrt{1+t} + \delta^2 n^{2\alpha}(1+t) \right) \\ &\leq n \left(D^* + n^{2\beta}\eta + \frac{1}{2} n^{-A-1} + \frac{1}{2} n^{-A-1} \right) = n(D^* + n^{2\beta}\eta) + n^{-A}. \end{aligned}$$

We denote $c', C_1, C_2, C_3, C_4, C_{12}$ to be explicit constants depending on $\varepsilon, \eta, R^*, \alpha$ (but not W or Σ_X). Now plug in $t = 2^{n\varepsilon} + C^{-1}n \ln 2(R^* + \varepsilon)$. The condition in Eq. (25) holds with probability at least

$$1 - \exp\left(-2^{n\varepsilon(1-cn^4(\alpha-\beta)^{-1})}\right) - \exp(n \ln 2(R^* + \varepsilon)) \exp(-Ct) = 1 - \exp\left(-2^{n\varepsilon(1-c'n^4(\alpha-\beta)^{-1})}\right).$$

Finally, the rate of this quantization scheme is $R = R^* + \varepsilon + \frac{1}{n} \log(1/\delta)$, which we now bound:

$$\begin{aligned} \log(1/\delta) &\leq \max \left\{ C_1 \log n + \frac{1}{2} \log(1+t) + 2\alpha \log n, C_2 \log n + \frac{1}{2} \log(1+t) + \alpha \log n \right\} \\ &\leq C_{12} \log n + \frac{1}{2} \log(1+t) + 2\alpha \log n. \end{aligned}$$

Plugging in the expression for t , we obtain

$$\log(1+t) \leq C_3 + n\varepsilon + C_4 \log n,$$

and therefore,

$$R \leq R^* + \frac{3}{2}\varepsilon + C \cdot \frac{\log n + \alpha \log n}{n}.$$

For sufficiently large n , the RHS is $\leq R^* + 2\varepsilon$. Moreover, for sufficiently large n , our final distortion bound simplifies to $n(D^* + n^{2\beta} \cdot 2\eta)$. Setting $\varepsilon = \varepsilon_o/2, \eta = \eta_o/2$, we conclude the proof. \blacksquare

Proof of Claim 1. Since for all $j \in [n]$, $\widetilde{W}_j^2 \leq \|\widetilde{W}\|_\infty^2$ and $\frac{\lambda_j}{1+\lambda_j T} \leq \frac{1}{T}$,

$$T \sum_j \frac{\widetilde{W}_j^2 \lambda_j^2}{(1 + \lambda_j T)^2} \leq \|\widetilde{W}\|_\infty^2 \sum_j \frac{\lambda_j}{1 + \lambda_j T},$$

yielding $0 \leq \tau \leq \|\widetilde{W}\|_\infty$. \square

Proof of Claim 2.

$$\begin{aligned}
 \sqrt{d_{\Sigma_X}(W, q(\tau)c_i)} &= \sqrt{\mathbb{E}_X \left[\|(W - q(\tau)c_i)^T X\|_2^2 \right]} \\
 &\leq \sqrt{d_{\Sigma_X}(W, \tau c_i)} + \sqrt{\mathbb{E}_X \left[\|(\delta_\tau \cdot c_i)^T X\|_2^2 \right]} \\
 &= \sqrt{d_{\Sigma_X}(W, \tau c_i)} + \delta_\tau \sqrt{c_i^T \Sigma_X c_i},
 \end{aligned}$$

where $\delta_\tau = |q(\tau) - \tau| \leq \delta n^\alpha$. By Hanson-Wright inequality (Prop. 6) for any codeword j , the second term can be bounded as

$$\Pr \left[|c_j^T \Sigma_X c_j - \text{tr} \Sigma_X| > nt \right] \leq 2 \exp \left(-c \min \left(\frac{t^2 n^2}{K^4 \|\Sigma_X\|_F^2}, \frac{tn}{K^2 \|\Sigma_X\|_{op}} \right) \right),$$

for universal constants c, K , and therefore, with probability at least $1 - \exp(-Ctn/\|\Sigma_X\|_{op}) \geq 1 - \exp(-Ct)$,

$$\delta_\tau \sqrt{c_j^T \Sigma_X c_j} \leq \delta n^\alpha \sqrt{n(1+t)}.$$

The statement of the Claim follows by a union bound over all $2^{n(R^*+\varepsilon)}$ codewords. \square

B.2. Proof of Lemma 9: Success Probability of Random Gaussian Code

Condition 1 (Admissible W, Σ_X) We consider $W \in \mathbb{R}^n, \Sigma_X \in \mathbb{S}_+^n$ satisfying

1. SVD $(\Sigma_X) = U \Lambda U^\top$ for $U \in \mathcal{O}_n$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ with $\text{tr}(\Lambda) = n$;
2. $\|U^\top W\|_\infty \leq n^\alpha$ for a known $\alpha > 0$.

For $W, C \in \mathbb{R}^n, \tau \in \mathbb{R}$, and $\Sigma \in \mathbb{S}_+^n$, recall that we define the distortion function to be

$$d_{\Sigma_X}(W, \tau C) = (W - \tau C)^\top \Sigma_X (W - \tau C). \quad (E \text{ def.})$$

For admissible (W, Σ_X) (Cond. 1), denote $T = T_{\text{rc}}^\lambda(R^*)$ to be a unique solution to $R_{\text{rc}}^\lambda(T) = R^*$ and define

$$\tau(W, \Sigma_X) = \left(T \sum_j \frac{(U^\top W)_j^2 \lambda_j^2}{(1 + \lambda_j T)^2} \right)^{1/2} \left(\sum_j \frac{\lambda_j}{1 + \lambda_j T} \right)^{-1/2}. \quad (\tau \text{ def.})$$

Lemma 9 (Gaussian Book Success) Let $\varepsilon, \eta > 0$ be constants and $\mathbf{C} = \{c_1, \dots, c_{M=2^{n(R^*+\varepsilon)}}\}$, $c_i \sim i.i.d. \mathcal{N}(0, I_n)$ be a randomly generated Gaussian codebook.

For admissible W, Σ_X (Cond. 1) with $\alpha < 1/4$,

$$\Pr_{\mathbf{C}} \left[\frac{1}{n} \min_{i \in [M]} d_{\Sigma_X}(W, \tau(W, \Sigma_X) \cdot c_i) \leq \mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*) + \eta \right] \geq 1 - \exp \left(-2^{n\varepsilon(1-cn^{4\alpha-1})} \right),$$

where $c = c(\varepsilon, \eta, R^*, \alpha)$ is an explicit constant function, $\mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*)$ is defined in Eq. (D_{rc}), and $\tau(W, \Sigma_X)$ is defined in Eq. (τ def.).

The proof of Lemma 9 relies on the following bound on a probability that a single codeword $c_i \sim \mathcal{N}(0, I_n)$ achieves small distortion d_{Σ_X} .

Lemma 10 *In the setting of Lemma 9, denote*

$$p_n := \Pr_{c_i \sim \mathcal{N}(0, I_n)} \left[\frac{1}{n} d_{\Sigma_X}(W, \tau(W, \Sigma_X) \cdot c_i) \leq \mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*) + \eta \right].$$

Then,

$$\ln p_n \geq -nR^*(1 + O(n^{4\alpha-1})) = -nR^*(1 + o(1)).$$

Proof (of Lemma 9) Denote $D^* = \mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*)$, and $\tau = \tau(W, \Sigma_X)$. Given the EVD $(\Sigma_X) = U\Lambda U^\top$, by definition of d_{Σ_X} , for all $i \in [M]$,

$$\begin{aligned} d_{\Sigma_X}(W, \tau \cdot c_i) &= (W - \tau c_i)^\top U\Lambda U^\top (W - \tau c_i) \\ &= (U^\top W - \tau U^\top c_i)^\top \Lambda (U^\top W - \tau U^\top c_i) \\ &=: (\widetilde{W} - \tau \tilde{c}_i)^\top \Lambda (\widetilde{W} - \tau \tilde{c}_i), \end{aligned}$$

where we denote $\widetilde{W} = U^\top W$ and $\tilde{c}_i = U^\top c_i$. Note that $\tilde{c}_i \sim \mathcal{N}(0, I_n)$ are independent.

We show that for the optimal choice of τ in Eq. ([\tau def.](#)) and $M = 2^{n(R^* + \varepsilon)}$, with probability at least $1 - \exp(-2^{n\varepsilon(1 - c\|\widetilde{W}\|_\infty^4/n)})$ (where $\|\widetilde{W}\|_\infty^4/n \leq n^{4\alpha-1}$),

$$\frac{1}{n} \min_{i \in [M]} d_{\Sigma_X}(W, \tau c_i) = \frac{1}{n} \min_{i \in [M]} \sum_j \lambda_j (\widetilde{W}_j - \tau \tilde{c}_{ij})^2 \leq D^* + \eta.$$

By Lemma 10, for any $i \in [M]$,

$$\ln p_n = \ln \Pr \left[\frac{1}{n} \sum_j \lambda_j (\widetilde{W}_j - \tau \tilde{c}_{ij})^2 \leq D^* + \eta \right] \geq -nR^*(1 + O(n^{4\alpha-1})).$$

Then, since c_1, \dots, c_M are independent,

$$\begin{aligned} \Pr \left[\frac{1}{n} \min_{i \in [M]} d_{\Sigma_X}(W, \tau c_i) \leq D^* + \eta \right] &\geq 1 - (1 - p_n)^M \geq 1 - \exp(-2^{-nR^*(1 + O(\|\widetilde{W}\|^4/n))} 2^{nR^* + n\varepsilon}) \\ &= 1 - \exp(-2^{n\varepsilon(1 - O(\|\widetilde{W}\|^4/n))}) \\ &\geq 1 - \exp(-2^{n\varepsilon(1 - O(n^{4\alpha-1}))}), \end{aligned}$$

which concludes the proof of Lemma 9. ■

Proof (of Lemma 10) Denote $R_{\text{nat}} = R^* \ln 2$, $D^* = \mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*)$, and $\tau = \tau(W, \Sigma_X)$. First, for all $j \in [n]$ denote $\mu_j = \mathcal{L} \left(\lambda_j (\widetilde{W}_j - \tau \tilde{c}_{ij})^2 \right)$, where $\tilde{c}_{ij} \sim \mathcal{N}(0, 1)$ and rewrite

$$p_n = \Pr_{X_j \sim \mu_j} \left[\frac{1}{n} \sum_j X_j \leq D^* + \eta \right].$$

Let $t = -T/(2\tau^2)$, where T, τ are defined in Eq. ([\tau def.](#)). We have

$$\ln \mathbb{E}_{X \sim \mu_j} [e^{tX}] = \frac{t\widetilde{W}_j^2 \lambda_j}{1 - 2t\lambda_j \tau^2} - \frac{1}{2} \ln(1 - 2t\lambda_j \tau^2) =: \Phi_j(t).$$

A direct computation yields

$$\Phi_j'(t) = \frac{\widetilde{W}_j^2 \lambda_j}{(1 - 2t\lambda_j \tau^2)^2} + \frac{\lambda_j \tau^2}{1 - 2t\lambda_j \tau^2}. \quad (27)$$

For each $j \in [m]$, define a new probability measure $\tilde{\mu}_j$ as $\frac{d\tilde{\mu}_j}{d\mu_j}(x) = e^{tx - \Phi_j(t)}$. We have $\int_{\mathbb{R}} d\tilde{\mu}_j = e^{-\Phi_j(t)} \int_{\mathbb{R}} e^{tx} d\mu_j = 1$. Since the exponential tilting above is applied to a noncentral χ^2 distribution μ , the transformation simply shifts/rescales the parameters of the χ^2 distribution. We can explicitly compute

$$\tilde{\mu}_j = \mathcal{L}\left(\mathcal{N}(m_j, s_j^2)\right), \quad \text{where } m_j = \frac{\widetilde{W}_j^2 \lambda_j^{1/2}}{1 - 2t\lambda_j \tau^2} \text{ and } s_j^2 = \frac{\lambda_j \tau^2}{1 - 2t\lambda_j \tau^2}.$$

Plugging in the choice of τ in Eq. ([\tau def.](#)) and $T = -2t\tau^2$ into Eq. (27), we obtain

$$\begin{aligned} \mathbb{E}_{X \sim \tilde{\mu}} \sum_j X_j &= \sum_j \int x e^{tx - \Phi_j(t)} d\mu_j(x) = \sum_j \frac{d}{dt} \ln \int e^{tx} d\mu_j(x) = \sum_j \Phi_j'(t) \\ &= \sum_j \frac{\widetilde{W}_j^2 \lambda_j}{1 - 2t\lambda_j \tau^2} = \sum_j \frac{\widetilde{W}_j^2 \lambda_j}{1 + \lambda_j T} = nD^*. \end{aligned}$$

Then,

$$D(\tilde{\mu} \parallel \mu) = \sum_j D(\tilde{\mu}_j \parallel \mu_j) = \sum_j t\Phi_j'(t) - \Phi_j = \frac{1}{2} \sum_j \ln(1 - 2t\lambda_j \tau^2) = \frac{1}{2} \sum_j \ln(1 + \lambda_j T) = nR_{\text{nat}}.$$

Recall that we defined $p_n = \Pr_{X_j \sim \mu_j} \left[\sum_j X_j \leq n(D^* + \eta) \right]$ and let $\tilde{p}_n = \Pr_{X_j \sim \tilde{\mu}_j} \left[\sum_j X_j \leq n(D^* + \eta) \right]$. By DPI,

$$nR_{\text{nat}} = D(\tilde{\mu} \parallel \mu) \geq d(\tilde{p}_n \parallel p_n) \geq -h(\tilde{p}_n) + \tilde{p}_n \ln \frac{1}{p_n} \geq -\ln 2 + \tilde{p}_n \ln \frac{1}{p_n},$$

which yields

$$\ln p_n \geq \frac{-nR_{\text{nat}} - \ln 2}{\tilde{p}_n}. \quad (28)$$

By Chebyshev's inequality,

$$1 - \tilde{p}_n = 1 - \Pr_{X_j \sim \tilde{\mu}_j} \left[\sum_j X_j \leq \mathbb{E} \sum_j X_j + n\eta \right] \leq \frac{\text{Var} \left[\sum_j X_j \right]}{n^2 \eta^2}.$$

In what follows we show that

$$\text{Var} \left[\sum_j X_j \right] = O \left(n \|\widetilde{W}\|_\infty^4 \right) = o(n^2), \quad (\text{Var})$$

yielding $1 - \tilde{p}_n = O \left(\|\widetilde{W}\|_\infty^4 / n \right) = o(1)$. Plugging this into Eq. (28), we get

$$\ln p_n \geq \frac{-nR_{\text{nat}} - \ln 2}{1 - O \left(\|\widetilde{W}\|_\infty^4 / n \right)} = -nR_{\text{nat}} (1 + O \left(\|\widetilde{W}\|_\infty^4 / n \right))$$

for sufficiently large n .

It remains to obtain the bound in Eq. (Var). From the definition of R_{nat} and $\frac{x}{1+x} \leq \ln(1+x) \leq x$ for $x > -1$,

$$|t|\tau^2 \sum_j \frac{\lambda_j}{1 - 2t\lambda_j\tau^2} \leq nR_{\text{nat}} = \frac{1}{2} \sum_j \ln(1 - 2t\lambda_j\tau^2) \leq |t|\tau^2 n.$$

Using $\frac{\lambda_j}{1 - 2t\lambda_j\tau^2} \leq \frac{1}{2|t|\tau^2}$ and $(1 - 2t\lambda_j\tau^2)^2 \geq 4|t|\lambda_j\tau^2$, we obtain

$$\begin{aligned} \text{Var} [X_j] &= \Lambda_j''(t) = \frac{4\widetilde{W}_j^2 \lambda_j^2 \tau^2}{(1 - 2t\lambda_j\tau^2)^3} + \frac{2\lambda_j^2 \tau^4}{(1 - 2t\lambda_j\tau^2)^2} \\ &\leq \frac{4\widetilde{W}_j^2 \lambda_j^2 \tau^2}{(1 - 2t\lambda_j\tau^2) \cdot 4|t|\lambda_j\tau^2} + \frac{\tau^2}{|t|} \cdot \frac{\lambda_j}{1 - 2t\lambda_j\tau^2} \\ &= \frac{1}{|t|} \cdot \frac{\widetilde{W}_j^2 \lambda_j}{1 - 2t\lambda_j\tau^2} + \frac{\tau^2}{|t|} \cdot \frac{\lambda_j}{1 - 2t\lambda_j\tau^2}, \end{aligned}$$

and therefore,

$$\text{Var} \left[\sum_j X_j \right] \leq n \cdot \left(\frac{D^*}{|t|} + \frac{R_{\text{nat}}}{|t|^2} \right). \quad (29)$$

We show a lower bound on parameter $|t|$. From the way we choose τ in Eq. (τ def.),

$$|t| = \frac{\sum_j \frac{\lambda_j}{1 - 2t\lambda_j\tau^2}}{\sum_j \frac{2\widetilde{W}_j^2 \lambda_j^2}{(1 - 2t\lambda_j\tau^2)^2}} \geq \frac{2R_{\text{nat}}}{\|\widetilde{W}\|_\infty^2} \frac{\sum_j \frac{\lambda_j \|\widetilde{W}\|_\infty^2}{1 - 2t\lambda_j\tau^2}}{\sum_j \frac{2\widetilde{W}_j^2 \lambda_j}{1 - 2t\lambda_j\tau^2}} \geq \frac{R_{\text{nat}}}{\|\widetilde{W}\|_\infty^2}, \quad (30)$$

where we used $\frac{\lambda_j}{1 - 2t\lambda_j\tau^2} \leq \frac{1}{2|t|\tau^2} \leq \frac{1}{2R_{\text{nat}}}$. Plugging (30) into the bound for variance in (29) and using $D^* \leq \frac{\|\widetilde{W}\|_\infty^2}{T} \leq \frac{\|\widetilde{W}\|_\infty^2}{2R_{\text{nat}}}$, we obtain

$$\text{Var} \left[\sum_j X_j \right] \leq n \cdot \|\widetilde{W}\|_\infty^4 \left(\frac{1}{2(R_{\text{nat}})^2} + \frac{1}{R_{\text{nat}}} \right) = O \left(n \|\widetilde{W}\|_\infty^4 \right),$$

which is the desired inequality in Eq. (Var). Here we used that $R_{\text{nat}} = \Theta(1)$. □

Appendix C. Random Coding: Gaussian Isotropic W

In this section we describe quantization of a Gaussian isotropic vector $W \sim \mathcal{N}(0, I_n) \in \mathbb{R}^n$ using random coding.

Rate-Distortion Function. Let $\lambda = (\lambda_1, \dots, \lambda_n)^\top \in \mathbb{R}^n$ be such that $\lambda_i \geq 0$ for all $i \in [n]$ and $\sum_i \lambda_i = n$. We define *random-coding rate-distortion function* in dimension n to be the following parametric curve for $T > 0$:

$$D_{\text{rc}}^\lambda(T) = \frac{1}{n} \sum_{i=1}^n \frac{\lambda_i}{1 + \lambda_i T} \quad \text{and} \quad R_{\text{rc}}^\lambda(T) = \frac{1}{2n} \sum_{i=1}^n \log(1 + \lambda_i T). \quad (\text{RDRC})$$

Throughout, \log denotes \log_2 and \ln denotes the natural logarithm. Denote $T_{\text{rc}}^\lambda(R)$ to be a unique value T , s.t. $R_{\text{rc}}^\lambda(T) = R$. Let

$$\mathbf{D}_{\text{rc}}^\lambda(R) \triangleq D_{\text{rc}}^\lambda(T_{\text{rc}}^\lambda(R)). \quad (\mathbf{D}_{\text{rc}})$$

We consider a task of quantizing a vector $W \sim \mathcal{N}(0, I_n) \in \mathbb{R}^n$ under a distortion function

$$d_{\Sigma_X}(W, \widehat{W}) = \mathbb{E}_X \left[\left\| \left(W - \widehat{W} \right)^T X \right\|_2^2 \right] = (W - \widehat{W})^\top \Sigma_X (W - \widehat{W}), \quad (d_{\Sigma_X})$$

where $X \in \mathbb{R}^n$ is a random vector with $\mathbb{E}XX^\top = \Sigma_X \in \mathbb{S}_+^n$. Here we assume EVD $(\Sigma_X) = U\Lambda U^\top$ for $U \in \mathcal{O}_n$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ with $\text{tr}(\Lambda) = n$. Our goal is to obtain an upper bound on $\mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W})]$.

In Theorem 11 we show that there exists an encoder-decoder that, for any constant target rate $R^* > 0$ and admissible Σ_X as above, achieves rate R^* and expected distortion $\mathbf{D}_{\text{rc}}^\lambda(R^*)$ asymptotically, with high probability over the randomness S shared between encoder and decoder (from which the codebook is generated).

Theorem 11 (Achievability of RC RD for Gaussian Input: Nonasymptotic Guarantee)

Fix any constants $R^*, \varepsilon_o, \eta_o, B > 0$. There exists an encoder $f : \mathbb{R}^n \times \mathbb{S}_+^n \times [0, 1] \rightarrow [2^{nR}]$, a decoder $g : [2^{nR}] \times [0, 1] \rightarrow \mathbb{R}^n$ with $R \leq R^* + \varepsilon_o$, and a (shared) random variable $S \in [0, 1]$ with the following property.

For $W \sim \mathcal{N}(0, I_n)$, any Σ_X (with EVD $(\Sigma_X) = U\Lambda U^\top$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$) as above we set $\widehat{W} = g(f(W, \Sigma_X, S), S)$. Then for any sufficiently large $n \geq n_0 = n_0(\varepsilon_o, \eta_o, R^*, B)$ we have

$$\Pr_S \left[\mathbb{E}_W \left[\frac{1}{n} d_{\Sigma_X}(W, \widehat{W}) \right] \leq \mathbf{D}_{\text{rc}}^\lambda(R^*) + \eta_o \right] \geq 1 - \exp(-n^B),$$

where $d_{\Sigma_X}(W, \widehat{W})$ is the distortion function in Eq. (d_{Σ_X}) .

In Theorem 3 we show that, in fact, the quantization scheme in Theorem 11 achieves the rate-distortion guarantee simultaneously for all $\Sigma_X \in \mathbb{S}_+^n$ with $\text{tr}(\Sigma_X) = n$, with high probability over the shared randomness (i.e., the codebook **C**)

Theorem 3 (Achievability of RC RD for Gaussian Input: Worst-Case Σ_X) Fix any constants $R^*, \varepsilon_o, \eta_o, B > 0$. There exist an encoder $f : \mathbb{R}^n \times \mathbb{S}_+^n \times [0, 1] \rightarrow [2^{nR}]$, a decoder $g : [2^{nR}] \times$

$[0, 1] \rightarrow \mathbb{R}^n$ with $R \leq R^* + \varepsilon_0$, and a (shared) random variable $S \in [0, 1]$ with the following property.

For $W \sim \mathcal{N}(0, I_n)$ and $\Sigma_X \in \mathbb{S}_+^n$ we set $\widehat{W}(\Sigma_X) = g(f(W, \Sigma_X, S), S)$. Then, for sufficiently large $n \geq n_0 = n_0(\varepsilon_0, \eta_0, R^*, B)$, we have

$$\Pr_S \left[\sup_{\substack{\Sigma_X \in \mathbb{S}_+^n \\ \text{tr}(\Sigma_X) = n}} \left(\frac{1}{n} \mathbb{E}_W \left[d_{\Sigma_X}(W, \widehat{W}(\Sigma_X)) \right] - \mathbf{D}_{\text{rc}}^{\text{spec}(\Sigma_X)}(R^*) \right) \leq \eta_0 \right] \geq 1 - \exp(-n^B),$$

where $d_{\Sigma_X}(W, \widehat{W})$ is the distortion function in Eq. (d $_{\Sigma_X}$).

C.1. Proof of Theorem 11

The proof uses the result of Theorem 7, which obtains a distortion guarantee for any fixed vector $W \in \mathbb{R}^n$. We adjust the quantization scheme to not rely on the Theorem 7 assumption $\|U^\top W\|_\infty \leq n^\alpha$ (Cond. 1): in the unlikely case of large $\|U^\top W\|_\infty$, the decoder returns 0.

Quantization Scheme. Fix any constant $\alpha \in (0, 1/4)$. From the shared randomness S, f, g generate a Gaussian codebook $\mathbf{C} = \{c_1, \dots, c_{M=2^n(R^*+\varepsilon)}\}$ for $\varepsilon = \varepsilon(\varepsilon_0)$ to be chosen later. Denote $\widetilde{W} = U^\top W$. We define:

- **Encoder f :** Let $T := T_{\text{rc}}^\lambda(R^*)$ and define the scaling parameter $\tau = \tau(\widetilde{W}, \lambda)$ as

$$\tau = \begin{cases} \left(T \sum_j \frac{\widetilde{W}_j^2 \lambda_j^2}{(1+\lambda_j T)^2} \right)^{1/2} \left(\sum_j \frac{\lambda_j}{1+\lambda_j T} \right)^{-1/2} & \text{if } \|\widetilde{W}\|_\infty \leq n^\alpha \\ 0 & \text{otherwise.} \end{cases} \quad (\tau \text{ def.})$$

We set f to be a tuple

$$f(W, \Sigma_X, S) = \left(\underset{i \in [M]}{\text{argmin}} d_{\Sigma_X}(W, \tau c_i), q(\tau) \right),$$

where $q(\tau) = \delta \|\widetilde{W}\|_\infty \lfloor \tau / (\delta \|\widetilde{W}\|_\infty) \rfloor$ is a rounding quantization scheme of precision δ and recall the distortion function

$$d_{\Sigma_X}(W, C) = (W - C)^\top \Sigma_X (W - C).$$

- **Decoder g :**

$$g(i, q(\tau), S) = q(\tau) \cdot c_i.$$

Rate-Distortion Bound. Denote In the quantization scheme above,

$$R = \underbrace{R^* + \varepsilon}_{\text{Gauss. codebook}} + \underbrace{\frac{1}{n} \log(1/\delta)}_{\tau \text{ quant.}}.$$

The rest of the proof is to obtain a high probability bound on the resulting distortion that, given $\widehat{W} = g(f(W, \Sigma_X, S), S) = q(\tau) \cdot c_i$, can be expressed as:

$$d_{\Sigma_X}(W, \widehat{W}) = d_{\Sigma_X}(W, q(\tau)c_i) = (W - q(\tau)c_i)^\top \Sigma_X (W - q(\tau)c_i).$$

In what follows, denote $D^* = \mathbf{D}_{\text{rc}}^\lambda(R^*)$. Let $\eta = \eta(\eta_o)$ be a constant to be chosen later. By Theorem 7, with an appropriate choice of parameter $\delta = \delta(\varepsilon, \eta, R^*, \alpha, n)$ in the quantization scheme above, we have for any fixed W such that $\|U^\top W\|_\infty \leq n^\alpha$ and sufficiently large $n \geq n_0 = n_0(\varepsilon, \eta, \alpha, R^*)$,

$$\Pr_{\mathbf{C}} \left[\frac{1}{n} d_{\Sigma_X}(W, \widehat{W}) \leq \mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*) + \eta \right] \geq 1 - \exp\left(-2^{n\varepsilon(1-cn^{4\alpha-1})}\right),$$

where $c = c(\varepsilon, \eta, \alpha, R^*)$. At the same time, the rate can be bounded as $R \leq R^* + \varepsilon_o$ (for an appropriately chosen $\varepsilon = \varepsilon(\varepsilon_o)$). Denote the event of the codebook failure as $F(W, \mathbf{C}) : \mathbb{R}^n \times (\mathbb{R}^n)^M \rightarrow \{0, 1\}$:

$$F(W, \mathbf{C}) = \left\{ \|U^\top W\|_\infty \leq n^\alpha \text{ and } \frac{1}{n} d_{\Sigma_X}(W, \widehat{W}) > \mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*) + \eta \right\}. \quad (F \text{ def.})$$

Rewriting the result of Theorem 7, we obtain that

$$\forall W \in \mathbb{R}^n : \|U^\top W\|_\infty \leq n^\alpha : \quad \Pr_{\mathbf{C}} [F(W, \mathbf{C})] \leq \exp\left(-2^{n\varepsilon(1-cn^{4\alpha-1})}\right).$$

Denote $E_{\mathbf{C}}$ to be the event $E_{\mathbf{C}} = \{\forall i \in [M] : \|c_i\|_2 \leq n^B\}$ for some fixed constant $B > 10$ and note that

$$\Pr_{\mathbf{C}} [E_{\mathbf{C}}^c] \leq M \cdot e^{-Cn^{2B}} = e^{n \ln 2(R^* + \varepsilon) - Cn^{2B}} \leq \exp(-c'n^{2B}),$$

for constant $c' = c'(\varepsilon, R^*)$ and sufficiently large (constant) n . We denote $\mathbf{C}|E_{\mathbf{C}}$ to be the distribution of \mathbf{C} conditioned on $E_{\mathbf{C}}$ and note that, since $\Pr_{\mathbf{C}} [E_{\mathbf{C}}] \geq 1/2$, Theorem 7 in fact yields

$$\forall W \in \mathbb{R}^n : \|U^\top W\|_\infty \leq n^\alpha : \quad \Pr_{\mathbf{C}|E_{\mathbf{C}}} [F(W, \mathbf{C})] \leq 2 \exp\left(-2^{n\varepsilon(1-c'n^{4\alpha-1})}\right).$$

We show in steps 1-2 that for any constant $A > 0$ and sufficiently large n ,

$$\Pr_{\mathbf{C}|E_{\mathbf{C}}} \left[\mathbb{E}_W \left[d_{\Sigma_X}(W, \widehat{W}) \right] \leq n(\mathbf{D}_{\text{rc}}^\lambda(R^*) + \eta) + n^{-A} \right] \geq 1 - \exp\left(-2^{n\varepsilon(1-c'n^{4\alpha-1})}\right),$$

and then conclude the proof in step 3 via a union bound.

Step 1: Bound on $\Pr_{\mathbf{C}|E_{\mathbf{C}}} [\Pr_{W \sim \mathcal{N}(0, I_n)} [F(W, \mathbf{C}) | \mathbf{C}] > p^*]$. Denote the event $E_W = \{\|U^\top W\|_\infty \leq n^\alpha\}$. All expressions below assume $W \sim \mathcal{N}(0, I_n)$.

$$\mathbb{E}_{\mathbf{C}|E_{\mathbf{C}}} \left[\Pr_W [F(W, \mathbf{C}) | \mathbf{C}] \right] = \Pr_{W, \mathbf{C}|E_{\mathbf{C}}} [F(W, \mathbf{C})] = \Pr_W \left[\Pr_{\mathbf{C}|E_{\mathbf{C}}} [F(W, \mathbf{C}) | W] \right] \leq \Pr_W [E_W] \cdot 2 \exp\left(-2^{n\varepsilon(1-cn^{4\alpha-1})}\right).$$

Then, by Markov's inequality,

$$\Pr_{\mathbf{C}|E_{\mathbf{C}}} \left[\Pr_W [F(W, \mathbf{C}) | \mathbf{C}] > p^* \right] \leq \frac{\Pr_W [E_W] \cdot 2 \exp\left(-2^{n\varepsilon(1-cn^{4\alpha-1})}\right)}{p^*}. \quad (31)$$

Step 2: Bound on $\mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W})]$ in case of codebook success in Eq. (31). Fix any codebook \mathbf{C} for which $\Pr_W [F(W, \mathbf{C}) | \mathbf{C}] \leq p^*$ and $E_{\mathbf{C}}$ hold. Expand

$$\mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W})] = \underbrace{\mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W}) \mathbb{1}_{E_W}]}_{\text{I}} + \underbrace{\mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W}) \mathbb{1}_{E_W^c}]}_{\text{II}}.$$

First,

$$\begin{aligned} \mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W}) \mathbb{1}_{E_W}] &\leq p^* \mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W}) \mathbb{1}_{E_W} | F(W, \mathbf{C})] + \\ &\quad \Pr_W [F(W, \mathbf{C})^c] \mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W}) \mathbb{1}_{E_W} | F(W, \mathbf{C})^c]. \end{aligned}$$

For any W , we have, since $q(\tau) \leq n^\alpha$ by Claim 1,

$$d_{\Sigma_X}(W, \widehat{W}) = (W - q(\tau)c_i)^\top \Sigma_X (W - q(\tau)c_i) \leq (\|U^\top W\|_\infty + n^\alpha \|U^\top c_i\|_\infty)^2 n,$$

and therefore, since the above is $\leq n^{3B}$ for $W \in E_W$ and $\mathbf{C} \in E_{\mathbf{C}}$,

$$p^* \mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W}) \mathbb{1}_{E_W} | F(W, \mathbf{C})] \leq p^* \cdot n^{3B}. \quad (\text{I.1})$$

Moreover, by the definition of $F(W, \mathbf{C})$ in Eq. (F def.),

$$\begin{aligned} \Pr_W [F(W, \mathbf{C})^c] \mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W}) \mathbb{1}_{E_W} | F(W, \mathbf{C})^c] &\leq n \mathbb{E}_W [\mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*) \cdot \mathbb{1}_{F(W, \mathbf{C})^c} \mathbb{1}_{E_W}] + n\eta \\ &\leq \mathbb{E}_W \left[\sum_i \frac{(U^\top W)_i^2 \lambda_i}{1 + T\lambda_i} \mathbb{1}_{E_W} \right] + n\eta \\ &\leq \mathbb{E}_W \left[\sum_i \frac{(U^\top W)_i^2 \lambda_i}{1 + T\lambda_i} \right] + n\eta \\ &= n(D^* + \eta). \end{aligned} \quad (\text{I.2})$$

In the above we used that $R_{\text{rc}}^\lambda(T)$ is independent of vector $U^\top W$, and therefore both $\mathbf{D}_{\text{rc}}^\lambda(U^\top W, R^*)$ and $\mathbf{D}_{\text{rc}}^\lambda(R^*)$ share the same parameter $T > 0$. Combining (I.1) and (I.2), we obtain

$$\text{I} \leq n(D^* + \eta) + p^* \cdot n^{3B}. \quad (\text{I})$$

Now in the case of E_W^c , we have $\tau = 0$, and therefore, $d_{\Sigma_X}(W, \widehat{W}) = W^\top \Sigma_X W \leq n \|U^\top W\|_\infty^2$. Note that $U^\top W \sim \mathcal{N}(0, I_n)$, so by the standard Gaussian tail bound,

$$\text{II} \leq n \mathbb{E}_W [\|U^\top W\|_\infty^2 \cdot \mathbb{1}_{\|U^\top W\|_\infty > n^\alpha}] = n \int_{t=n^\alpha}^\infty t^2 \cdot 2ne^{-t^2/2} \leq 4n^{2+\alpha} e^{-n^{2\alpha}/2}. \quad (\text{II})$$

Combining (I) and (II) together, we obtain

$$\mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W})] \leq n(D^* + \eta) + p^* \cdot n^{3B} + 4n^{2+\alpha} e^{-n^{2\alpha}/2} \leq n(D^* + \eta) + n^{-A},$$

for sufficiently large constant n and the choice $p^* = e^{-n}$.

Step 3: Bound on the probability of successful codebook \mathbf{C} . Recall that $\Pr_{\mathbf{C}} [E_{\mathbf{C}}^c] \leq \exp(-c'n^{2B})$. From Eq. (31) and the choice of p^* above, we have for sufficiently large constant n ,

$$\Pr_{\mathbf{C}|E_{\mathbf{C}}} \left[\Pr_W [F(W, \mathbf{C}) | \mathbf{C}] > p^* \right] \leq \exp\left(-2^{n\varepsilon(1-cn^{4\alpha-1})} + n\right) \leq \exp(-c'n^{2B}).$$

Then, by a union bound, the probability to select a good codebook is at least $1 - \exp(-c'n^{2B})$, so we conclude (since for sufficiently large n , $n^{-A} < n\eta$),

$$\Pr_{\mathbf{C}} \left[\mathbb{E}_W \left[\frac{1}{n} d_{\Sigma_X}(W, \widehat{W}) \right] \leq D^* + 2\eta \right] \geq 1 - \exp(-c'n^{2B}),$$

from which the statement of the theorem follows.

C.2. Proof of Theorem 3

We show Theorem 3 by applying the results of Theorem 11 for fixed Σ_X with a covering argument. The quantization scheme is the same as in the proof of Theorem 11; the technical challenge of Theorem 3 is to show that it succeeds simultaneously with high probability for all $\Sigma_X \in \mathbb{S}_+^n$ with $\text{tr}(\Sigma_X) = n$.

We first recall some notation. In the quantization scheme of Theorem 11 we have a Gaussian codebook $\mathbf{C} = \{c_1, \dots, c_M\}$ for $M = 2^{n(R^* + \varepsilon)}$ (where $\varepsilon = \varepsilon(\varepsilon_0)$ is a constant) that is generated from the shared randomness $S \in [0, 1]$. We denote $\text{EVD}(\Sigma_X) = U\Lambda U^\top$ for $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$, and recall the distortion

$$d_{\Sigma_X}(W, C) = (W - C)^\top \Sigma_X (W - C).$$

The quantizer then sets the optimal scaling factor $\tau = \tau(U^\top W, \Lambda)$ in Eq. ([\tau def.](#)) and returns a tuple

$$f(W, \Sigma_X, S) = \left(\underset{i \in [M]}{\text{argmin}} d_{\Sigma_X}(W, \tau c_i), q(\tau) \right),$$

where $q(\tau)$ is a simple rounding quantizer. The decoder is defined to be

$$g((i, q(\tau)), S) = q(\tau) \cdot c_i.$$

Finally, for some constant $B > 10$, we denote

$$E_{\mathbf{C}} = \{\forall i \in [M] : \|c_i\|_2 \leq n^B\}, \quad \Pr_{\mathbf{C}} [E_{\mathbf{C}}^c] \leq \exp(-c'n^{2B})$$

and, denoting the conditional distribution $\mathbf{C}|E_{\mathbf{C}}$, the proof of Theorem 11 shows for any $\alpha \in (0, 1/4)$ and $\eta = \eta(\eta_0)$ to be chosen later,

$$\Pr_{\mathbf{C}|E_{\mathbf{C}}} \left[\mathbb{E}_W \left[\frac{1}{n} d_{\Sigma_X}(W, \widehat{W}) \right] \leq \mathbf{D}_{\text{rc}}^\lambda(R^*) + \eta \right] \geq 1 - \exp\left(-2^{n\varepsilon(1-c'n^{4\alpha-1})}\right), \quad (32)$$

for sufficiently large $n \geq n_0 = n_0(\varepsilon, \eta, R^*, \alpha)$. Simultaneously, it is shown that the scheme above achieves final rate $R \leq R + \varepsilon_0$ for the appropriately chosen $\varepsilon = \varepsilon(\varepsilon_0)$.

Step 1: Covering for $U \in \mathcal{O}_n$ and Λ .

Fact 1 (Covering of \mathcal{O}_n) For a universal constant c ,

$$N(\gamma, \mathcal{O}_n, \|\cdot\|_{op}) \leq (c/\gamma)^{n^2}.$$

Consequently, there exists a γ -covering $U_1, \dots, U_N \in \mathcal{O}_n$ such that for all $U \in \mathcal{O}_n$,

$$\min_{i \in N} \|U_i - U\|_{op} \leq \gamma,$$

and $N \leq (c/\gamma)^{n^2}$.

Fact 2 (Covering of Λ) For a universal constant c ,

$$N\left(\delta/\sqrt{d}; B_1^d, \|\cdot\|_2\right) \leq (c + c/\delta)^d.$$

Consequently, there exists a $(\tilde{\gamma}\sqrt{n})$ -covering $\tilde{s}_1, \dots, \tilde{s}_{N'} \in \left\{s \in \mathbb{R}_+^n : \sum_j s_j \leq n\right\}$ such that for all $\lambda \in \mathbb{R}_+^n$ with $\sum_j \lambda_j \leq n$,

$$\min_{i \in N'} \|\tilde{s}_i - \lambda\|_2 \leq \tilde{\gamma}\sqrt{n},$$

and $N' \leq (c/\tilde{\gamma})^n$ for a universal constant c .

Setting $s_i := \operatorname{argmin}_{s \in \mathbb{R}_+^n : \sum_j s_j = n} \|s - \tilde{s}_i\|_2$ and applying Fact 2 with $\tilde{\gamma} = \gamma/2$, we obtain:

Corollary 12 (Of Fact 2) There exists a $(\gamma\sqrt{n})$ -covering $s_1, \dots, s_{N'} \in \left\{s \in \mathbb{R}_+^n : \sum_j s_j = n\right\}$ such that for all $\lambda \in \mathbb{R}_+^n$ with $\sum_j \lambda_j = n$,

$$\min_{i \in N'} \|s_i - \lambda\|_2 = \gamma\sqrt{n},$$

and $N' \leq (c/\gamma)^n$ for a universal constant c .

Consider the γ -covering $U_1, \dots, U_N \in \mathcal{O}_n$ and $(\gamma_2\sqrt{n})$ -covering $s_1, \dots, s_{N'} \in \left\{s \in \mathbb{R}_+^n : \sum_j s_j = n\right\}$ from the facts above for γ to be chosen later. Denote $\Lambda_i := \operatorname{diag}(s_i)$. By a union bound applied to Eq. (32),

$$\Pr_{\mathbf{C}|E_C} \left[\sup_{\Sigma_X = U_i \Lambda_j U_i^\top} \left(\mathbb{E}_W \left[\frac{1}{n} d_{\Sigma_X}(W, \widehat{W}) \right] - \mathbf{D}_{\text{rc}}^{s_j}(R^\star) \right) \leq \eta \right] \geq 1 - NN' \exp\left(-2^{n\varepsilon(1-c'n^{4\alpha-1})}\right).$$

Step 2: Perturbation bound on $\mathbb{E}_W [d_{\Sigma_X}(W, \widehat{W})]$. Fix any constant $A > 0$ (say $A = 10$). We show the following distortion perturbation results.

Claim 3 (Distortion Perturbation) Let $\Lambda = \operatorname{diag}(\lambda)$, $\tilde{\Lambda} = \operatorname{diag}(\tilde{\lambda})$ be s.t. $\Lambda, \tilde{\Lambda} \succeq 0$ and $\operatorname{tr}(\Lambda) = \operatorname{tr}(\tilde{\Lambda}) = n$ and $U, \tilde{U} \in \mathcal{O}_n$ satisfy

$$\|\Lambda - \tilde{\Lambda}\|_2 \leq \gamma\sqrt{n} \quad \text{and} \quad \|U - \tilde{U}\|_{op} \leq \gamma.$$

Denote

$$\Sigma_X = U\Lambda U^\top \quad \text{and} \quad \tilde{\Sigma}_X = \tilde{U}\tilde{\Lambda}\tilde{U}^\top.$$

If $\mathbf{C} = \{c_1, \dots, c_M\} \in (\mathbb{R}^n)^M$ satisfies $\max_{i \in [M]} \|c_i\|_2 \leq n^B$, then, in the setup above,

$$\mathbb{E}_{W \sim \mathcal{N}(0, I_n)} \left[d_{\tilde{\Sigma}_X}(W, \widehat{W}) \right] \leq \mathbb{E}_{W \sim \mathcal{N}(0, I_n)} \left[d_{\Sigma_X}(W, \widehat{W}) \right] + \frac{1}{2}n^{-A} + \gamma n^{O(B)} \cdot 2^{4nR^*},$$

for any sufficiently large (constant) n , where $O(B)$ hides constants in A, B, R^* .

Plugging in $\gamma = 2^{-5nR^*}$, the bound in Claim 3 simplifies to n^{-A} , since for sufficiently large (constant) n , $\gamma n^{O(B)} \cdot 2^{4nR^*} \leq \frac{1}{2}n^{-A}$. Moreover, as we will shortly see, $|\mathbf{D}_{\text{rc}}^\lambda(R^*) - \mathbf{D}_{\text{rc}}^{s_j}(R^*)| \leq \text{poly}(n) \cdot \|\lambda - s_j\|_2 \leq \text{poly}(n) \cdot \gamma$, so, for a fixed codebook \mathbf{C} and sufficiently large n ,

$$\sup_{\Sigma_X = U\Lambda U^\top, U \in \mathcal{O}_n} \left(\mathbb{E}_W \left[d_{\Sigma_X}(W, \widehat{W}) \right] - n\mathbf{D}_{\text{rc}}^\lambda(R^*) \right) \leq \sup_{\Sigma_X = U_i \Lambda_j U_i^\top} \left(\mathbb{E}_W \left[d_{\Sigma_X}(W, \widehat{W}) \right] - n\mathbf{D}_{\text{rc}}^{s_j}(R^*) \right) + n^{-A}.$$

Then,

$$\Pr_{\mathbf{C} | E_{\mathbf{C}}} \left[\sup_{\Sigma_X = U\Lambda U^\top, U \in \mathcal{O}_n} \left(\mathbb{E}_W \left[\frac{1}{n} d_{\Sigma_X}(W, \widehat{W}) \right] - \mathbf{D}_{\text{rc}}^\lambda(R^*) \right) \leq \eta + n^{-A-1} \right] \geq 1 - NN' \exp \left(-2^{n\varepsilon(1-c'n^{4\alpha-1})} \right).$$

In the above, $NN' \leq (c/\gamma)^{n^2+n} = 2^{5(n^3+n^2)R^*}$, so the probability on the RHS of the equation above is $\geq 1 - \exp \left(-2^{n\varepsilon(1-cn^{4\alpha-1})} \right)$ for a constant c . The statement of Theorem 3 follows by a union bound and the fact that $\Pr_{\mathbf{C}} [E_{\mathbf{C}}^c] \leq \exp(-c'n^{2B})$ and choosing n such that $n^{-A-1} \leq \eta$.

It remains to verify that $|\mathbf{D}_{\text{rc}}^\lambda(R^*) - \mathbf{D}_{\text{rc}}^{s_j}(R^*)| \leq \text{poly}(n) \cdot \|\lambda - s_j\|_2$. An explicit calculation (using the implicit function theorem) gives for all $j \in [n]$

$$\begin{aligned} \frac{\partial n\mathbf{D}_{\text{rc}}^\lambda(R^*)}{\partial \lambda_j} &= \frac{1}{(1 + \lambda_j T)^2} + \sum_i \frac{\lambda_i^2}{(1 + \lambda_i T)^2} \cdot \frac{\frac{T}{1 + \lambda_j T}}{\sum_i \frac{\lambda_i}{1 + \lambda_i T}} \\ &\leq \frac{1}{(1 + \lambda_j T)^2} + \frac{1}{1 + \lambda_j T} \leq 2, \end{aligned}$$

which yields the desired statement.

Proof (of Claim 3) Denote $\tau = \tau(W, \Sigma_X)$ and $\tilde{\tau} = \tau(W, \tilde{\Sigma}_X)$. In our notation the claim statement is equivalent to, for $W \sim \mathcal{N}(0, I_n)$,

$$\mathbb{E}_W \left[d_{\tilde{\Sigma}_X}(W, q(\tilde{\tau})c_i) \right] \leq \mathbb{E}_W \left[d_{\Sigma_X}(W, q(\tau)c_j) \right] + \frac{1}{2}n^{-A} + \gamma n^{O(B)} \cdot 2^{4nR^*},$$

where $i = \text{argmin}_k d_{\tilde{\Sigma}_X}(W, \tilde{\tau}c_k)$ and $j = \text{argmin}_k d_{\Sigma_X}(W, \tau c_k)$.

Error from quantizing $\tau, \tilde{\tau}$. We first uniformly bound the effect of quantizing $\tau, \tilde{\tau}$. Similarly to Claim 2,

$$\left| \sqrt{d_{\Sigma_X}(W, q(\tau)c_j)} - \sqrt{d_{\Sigma_X}(W, \tau c_j)} \right| \leq \delta_\tau \sqrt{c_j^\top \Sigma_X c_j} \leq \delta_\tau n^{B+1}.$$

A standard Gaussian tail argument, combined with a bound $\max\{q(\tau), \tau\} \leq \|U^\top W\|_\infty$, shows that $\mathbb{E}_W \left[\sqrt{d_{\Sigma_X}(W, q(\tau)c_j)} + \sqrt{d_{\Sigma_X}(W, \tau c_j)} \right] \leq n^P$ for a sufficiently large constant P , so setting $\delta_\tau \leq \frac{1}{8}n^{-B-1-P-A}$ incurs $o_n(1)$ factors in R (see proof of Theorem 7) and achieves

$$\left| d_{\Sigma_X}(W, q(\tau)c_j) - d_{\Sigma_X}(W, \tau c_j) \right| \leq \frac{1}{8}n^{-A} \quad \text{and} \quad \left| d_{\tilde{\Sigma}_X}(W, q(\tilde{\tau})c_i) - d_{\tilde{\Sigma}_X}(W, \tilde{\tau}c_i) \right| \leq \frac{1}{8}n^{-A}.$$

Thus, it is sufficient to show

$$\mathbb{E}_W \left[\min_i d_{\tilde{\Sigma}_X}(W, \tilde{\tau}c_i) \right] \leq \mathbb{E}_W \left[\min_i d_{\Sigma_X}(W, \tau c_i) \right] + \frac{1}{4}n^{-A} + \gamma n^{O(B)} \cdot 2^{4nR^*}.$$

Error from the tail W event. First, denote the event $E_W = \{\|W\|_2 \leq n^B\}$. We can expand

$$\mathbb{E}_W \left[\min_i d_{\Sigma_X}(W, \tau c_i) \right] = \mathbb{E}_W \left[\min_i d_{\Sigma_X}(W, \tau c_i) \cdot \mathbb{1}_{E_W} \right] + \mathbb{E}_W \left[\min_i d_{\Sigma_X}(W, \tau c_i) \cdot \mathbb{1}_{E_W^c} \right],$$

and, by definition of $d_{\Sigma_X}(W, \tau c_i)$, the second term is

$$\begin{aligned} \mathbb{E}_W \left[\min_i d_{\Sigma_X}(W, \tau c_i) \cdot \mathbb{1}_{E_W^c} \right] &= \mathbb{E}_W \left[\min_i (W - \tau c_i)^\top \Sigma_X (W - \tau c_i) \cdot \mathbb{1}_{E_W^c} \right] \\ &\leq \mathbb{E}_W \left[n \|U^\top W\|_\infty^2 (1 + \|U^\top c_i\|_\infty) \cdot \mathbb{1}_{E_W^c} \right] \\ &\leq 2n^{21} \mathbb{E}_W \left[\|U^\top W\|_\infty^2 \cdot \mathbb{1}_{E_W^c} \right] \\ &\leq 2n^{2B+1} \mathbb{E}_W \left[\|W\|_2^2 \cdot \mathbb{1}_{E_W^c} \right] \\ &\leq 2n^{2B+1} \int_{t=n^B}^\infty t e^{-c_\chi t} \leq c' n^{3B} e^{-c' n^B} \leq \frac{1}{4} n^{-A} \end{aligned}$$

for sufficiently large n , where we used $\tau \leq \|U^\top W\|_\infty$, see Claim 1, and c_χ, c' are universal constants.

Error from Σ_X perturbation. We now show that

$$\mathbb{E}_W \left[\min_i d_{\tilde{\Sigma}_X}(W, \tilde{\tau}c_i) \cdot \mathbb{1}_{E_W} \right] \leq \mathbb{E}_W \left[\min_i d_{\Sigma_X}(W, \tau c_i) \cdot \mathbb{1}_{E_W} \right] + \gamma n^{O(B)} \cdot 2^{4nR^*}, \quad (33)$$

which, combined with the bound above yields the Claim. We use Lemma 14 (proved in Sec. E), which bounds $|\tau(W, \Sigma_X) - \tau(W, \tilde{\Sigma}_X)|$ in terms of $\|U^\top W - \tilde{U}^\top W\|_\infty$ and $\|\Lambda - \tilde{\Lambda}\|_2$. Given that $\|U^\top W\|_\infty, \|\tilde{U}^\top W\|_\infty \leq \|W\|_2 \leq n^B$ under the E_W event, Lemma 14 gives a (crude) bound:

$$|\tau - \tilde{\tau}| = |\tau(W, \Sigma_X) - \tau(W, \tilde{\Sigma}_X)| \leq \|U^\top W - \tilde{U}^\top W\|_\infty + n^{O(B)} \|\Lambda - \tilde{\Lambda}\|_2 \cdot 2^{4nR^*}.$$

Plugging in $\|U^\top W - \tilde{U}^\top W\|_\infty \leq \|U - \tilde{U}\|_{op} \|W\|_2 \leq \gamma n^B$ and $\|\Lambda - \tilde{\Lambda}\|_2 \leq \gamma \sqrt{n}$, we obtain

$$|\tau - \tilde{\tau}| \leq \gamma n^{O(B)} \cdot 2^{4nR^*}.$$

Additionally,

$$\begin{aligned} \|\tilde{\Sigma}_X - \Sigma_X\|_{op} &\leq \|U^\top (\Lambda - \tilde{\Lambda}) U\|_{op} + \|U^\top \tilde{\Lambda} U - \tilde{U}^\top \tilde{\Lambda} \tilde{U}\|_{op} \\ &\leq \|\Lambda - \tilde{\Lambda}\|_2 + \|(U - \tilde{U})^\top \tilde{\Lambda} \tilde{U}\|_{op} + \|\tilde{U}^\top \tilde{\Lambda} (U - \tilde{U})\|_{op} + \|(U - \tilde{U})^\top \tilde{\Lambda} (U - \tilde{U})\|_{op} \\ &\leq \gamma \sqrt{n} + 2\gamma n + \gamma^2 n \leq 4\gamma n. \end{aligned}$$

Again denote $j = \operatorname{argmin}_k d_{\Sigma_X}(W, \tau c_k)$, and let $x = W - \tau c_j, \tilde{x} = W - \tilde{\tau} c_j$. Then,

$$\begin{aligned} \min_i d_{\tilde{\Sigma}_X}(W, \tilde{\tau}c_i) \cdot \mathbb{1}_{E_W} &\leq d_{\tilde{\Sigma}_X}(W, \tilde{\tau}c_j) \cdot \mathbb{1}_{E_W} \\ &= (W - \tilde{\tau}c_j)^\top \tilde{\Sigma}_X (W - \tilde{\tau}c_j) \cdot \mathbb{1}_{E_W} \\ &= \tilde{x}^\top \tilde{\Sigma}_X \tilde{x} \cdot \mathbb{1}_{E_W} \\ &= x^\top \Sigma_X x \cdot \mathbb{1}_{E_W} + (\tilde{x} - x)^\top \tilde{\Sigma}_X \tilde{x} \cdot \mathbb{1}_{E_W} + x^\top \Sigma_X (\tilde{x} - x) \cdot \mathbb{1}_{E_W} \\ &\quad + x^\top (\tilde{\Sigma}_X - \Sigma_X) \tilde{x} \cdot \mathbb{1}_{E_W}. \end{aligned}$$

Since $\max\{\tau, \tilde{\tau}\} \leq \max\{\|U^\top W\|_\infty, \|\tilde{U}^\top W\|_\infty\} \leq \|W\|_2$, we can bound

$$\|x\|_2 \cdot \mathbb{1}_{E_W}, \|\tilde{x}\|_2 \cdot \mathbb{1}_{E_W} \leq \|W\|_2(1 + \|C_j\|_2) \cdot \mathbb{1}_{E_W} \leq 2n^{2B}$$

and, using the bound on $|\tau - \tilde{\tau}|$ above,

$$\|x - \tilde{x}\|_2 \leq |\tau - \tilde{\tau}| \cdot n^B \leq \gamma n^{O(B)} \cdot 2^{4nR^*}.$$

Plugging this into our bound, we obtain

$$\begin{aligned} \min_i d_{\tilde{\Sigma}_X}(W, \tilde{\tau}c_i) \cdot \mathbb{1}_{E_W} &\leq \min_i d_{\Sigma_X}(W, \tau c_i) \cdot \mathbb{1}_{E_W} + 2\gamma n^{O(B)} \cdot 2^{4nR^*} + 4n^{4B} \cdot 4\gamma n \\ &\leq \min_i d_{\Sigma_X}(W, \tau c_i) \cdot \mathbb{1}_{E_W} + \gamma n^{O(B)} \cdot 2^{4nR^*}, \end{aligned}$$

which concludes the proof of this step and the claim. ■

□

Appendix D. Worst-Case Gap Between Waterfilling and Random Coding: Theorem 5

Recall that for $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ with $\text{tr}(\Lambda) = n$, the random-coding rate-distortion is given by a parametric relationship

$$D_{\text{rc}}(\Lambda, T) = \frac{1}{n} \sum_{i=1}^n \frac{\lambda_i}{1 + \lambda_i T} \quad R_{\text{rc}}(\Lambda, T) = \frac{1}{2n} \sum_{i=1}^n \log(1 + \lambda_i T), \quad (\text{RDRC})$$

and the waterfilling rate-distortion by

$$D_{\text{wf}}(\Lambda, t) = \frac{1}{n} \sum_{i=1}^n \min\{\lambda_i, t\} \quad R_{\text{wf}}(\Lambda, t) = \frac{1}{2n} \sum_{i=1}^n \max\{0, \log(\lambda_i/t)\}. \quad (\text{RDWF})$$

Throughout, \log denotes the base-2 logarithm, and \ln denotes the natural logarithm. The goal of this section is to quantify, at a fixed target distortion D^* , the rate overhead of our universal codebook, whose rate-distortion is given in Eq. (RDRC), compared to the Σ_X -fine-tuned optimal codebook, whose rate-distortion is given in Eq. (RDWF).

Denote⁵ the implicit functions $T_{\text{rc}}(\Lambda, D^*) = T^*$ s.t. $D_{\text{rc}}(\Lambda, T^*) = D^*$ and $t_{\text{wf}}(\Lambda, D^*) = t^*$ s.t. $D_{\text{wf}}(\Lambda, t^*) = D^*$. Denote

$$\mathbf{R}_{\text{rc}}(\Lambda, D^*) = R_{\text{rc}}(\Lambda, T_{\text{rc}}(\Lambda, D^*)) \quad \text{and} \quad \mathbf{R}_{\text{wf}}(\Lambda, D^*) = R_{\text{wf}}(\Lambda, t_{\text{wf}}(\Lambda, D^*)).$$

In this notation, the rate overhead incurred by our universal codebook is

$$\sup_{\Lambda} \{\mathbf{R}_{\text{rc}}(\Lambda, D^*) - \mathbf{R}_{\text{wf}}(\Lambda, D^*)\},$$

where the supremum is over $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ with $\text{tr}(\Lambda) = n$. In Theorem 5, we prove

$$\sup_{D^* \in (0,1)} \sup_{\Lambda} \{\mathbf{R}_{\text{rc}}(\Lambda, D^*) - \mathbf{R}_{\text{wf}}(\Lambda, D^*)\} \leq 0.11.$$

Concretely, we first verify that the maximum rate gap is approached at spectra containing at most 2 distinct eigenvalue and at vanishing distortions. The resulting gap expression can be directly bounded by 0.11 bit.

5. All of the $D_{\text{rc}}, D_{\text{wf}}, R_{\text{rc}}, R_{\text{wf}}$ are monotone in T, t .

D.1. Proof of Theorem 5
Theorem 5 (Worst-Case Gap Between Waterfilling and Random Coding)

$$\sup_{D^* \in (0,1)} \sup_{\Lambda} \{\mathbf{R}_{\text{rc}}(\Lambda, D^*) - \mathbf{R}_{\text{wf}}(\Lambda, D^*)\} \leq 0.11,$$

where the supremum is over $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ with $\text{tr}(\Lambda) = n$.

Theorem 5 follows immediately from a more general version of the same statement in Theorem 13. Let $\mathcal{P}([0, \infty))$ denote the space of Borel probability measures on $[0, \infty)$, and define

$$\mathcal{P}_1([0, \infty)) = \left\{ \lambda \in \mathcal{P}([0, \infty)) : \int x d\lambda(x) = 1 \right\}.$$

We extend the finite-dimension rate-distortion curves in Eq. (RDRC) and (RDWF) as follows: for $\mu \in \mathcal{P}_1([0, \infty))$, let

$$D_{\text{rc}}(\mu, T) = \int_0^\infty \frac{x}{1+xT} d\mu(x), \quad R_{\text{rc}}(\mu, T) = \frac{1}{2} \int_0^\infty \log(1+xT) d\mu(x),$$

and

$$D_{\text{wf}}(\mu, t) = \int_0^\infty \min\{x, t\} d\mu(x), \quad R_{\text{wf}}(\mu, t) = \frac{1}{2} \int_0^\infty \max\{0, \log(x/t)\} d\mu(x).$$

Similarly to above, for each $D^* \in (0, 1)$, let $T_{\text{rc}}(\mu, D^*)$ and $t_{\text{wf}}(\mu, D^*)$ be defined by $D_{\text{rc}}(\mu, T_{\text{rc}}(\mu, D^*)) = D^*$ and $D_{\text{wf}}(\mu, t_{\text{wf}}(\mu, D^*)) = D^*$, and set

$$\mathbf{R}_{\text{rc}}(\mu, D^*) = R_{\text{rc}}(\mu, T_{\text{rc}}(\mu, D^*)), \quad \mathbf{R}_{\text{wf}}(\mu, D^*) = R_{\text{wf}}(\mu, t_{\text{wf}}(\mu, D^*)).$$

Theorem 13 (Worst-Case Gap Between Waterfilling and Random Coding)

$$\sup_{D^* \in (0,1)} \sup_{\mu} \{\mathbf{R}_{\text{rc}}(\mu, D^*) - \mathbf{R}_{\text{wf}}(\mu, D^*)\} \leq 0.11,$$

where the supremum is over probability measures $\mu \in \mathcal{P}_1([0, \infty))$.

Indeed, for any $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \succeq 0$ with $\text{tr}(\Lambda) = n$, the empirical measure $\mu_\Lambda = \frac{1}{n} \sum_{i=1}^n \delta_{\lambda_i}$ belongs to $\mathcal{P}_1([0, \infty))$.

Proof (of Theorem 13). Denote the rate-gap

$$\Delta(\mu, D^*) = \frac{1}{2} \int r_{T_{\text{rc}}(\mu, D^*), t_{\text{wf}}(\mu, D^*)}(x) d\mu(x), \quad r_{T,t}(x) = \begin{cases} \log(1+xT) & \text{if } x \leq t \\ \log\left(\frac{1+xT}{x/t}\right) & \text{if } x > t. \end{cases}$$

It will be convenient to consider the following rescaling:

$$\tilde{\mu} := (x \rightarrow x/D^*)_{\#}\mu, \quad (\tilde{\mu})$$

so that $\int x d\tilde{\mu}(x) = (D^*)^{-1} > 1$. Accordingly, we define $D^*(\tilde{\mu}) = (\int x d\tilde{\mu}(x))^{-1}$. Additionally, rescale the implicitly defined parameters t_{wf} and T_{rc} :

$$\tilde{t}_{\text{wf}}(\tilde{\mu}) \triangleq t_{\text{wf}}((x \rightarrow x \cdot D^*)_{\#\tilde{\mu}}, D^*(\tilde{\mu})) / D^*(\tilde{\mu}) \quad \text{and} \quad \tilde{T}_{\text{rc}}(\tilde{\mu}) \triangleq T_{\text{rc}}((x \rightarrow x \cdot D^*)_{\#\tilde{\mu}}, D^*(\tilde{\mu})) \cdot D^*(\tilde{\mu}),$$

i.e., $\tilde{t}_{\text{wf}}(\tilde{\mu}) = t_{\text{wf}}(\mu, D^*)/D^*$ and $\tilde{T}_{\text{rc}}(\tilde{\mu}) = T_{\text{rc}}(\mu, D^*) \cdot D^*$ in the original parameters. The distortion conditions then correspond to

$$1 = \int \frac{x}{1+x\tilde{T}} d\tilde{\mu}(x) = \int \min\{x, \tilde{t}\} d\tilde{\mu}(x), \quad \text{where } \tilde{t} = \tilde{t}_{\text{wf}}(\tilde{\mu}), \tilde{T} = \tilde{T}_{\text{rc}}(\tilde{\mu}). \quad (\text{Dist})$$

We can now rewrite the optimization objective in terms of $\tilde{\mu}$:

$$\Phi(\tilde{\mu}) \triangleq \Delta((x \rightarrow x \cdot D^*)_{\# \tilde{\mu}}, D^*(\tilde{\mu})) = \Delta(\mu, D^*) = \frac{1}{2} \int r_{\tilde{T}, \tilde{t}}(x) d\tilde{\mu}(x). \quad (\Phi)$$

Finally, for technical reasons, we will consider $\tilde{\mu} \in \mathcal{P}([0, \infty]) \supseteq \mathcal{P}([0, \infty])$ with the natural extensions of functions $\frac{x}{1+x\tilde{T}}, \log\left(\frac{1+x\tilde{T}}{x/\tilde{t}}\right)$. To prove the theorem it is then sufficient to show

$$\sup_{\tilde{\mu} \in \mathcal{A}} \Phi(\tilde{\mu}) \leq 0.11, \quad \text{where } \mathcal{A} \triangleq \left\{ \tilde{\mu} \in \mathcal{P}([0, \infty]) : D^*(\tilde{\mu}) \in (0, 1) \right\}.$$

In what follows, we will restrict the optimization scope \mathcal{A} in three steps. For $\tilde{\mu}$ as above, to simplify notation, when clear from the context, we write $D^* = D^*(\tilde{\mu})$, $\tilde{t} = \tilde{t}_{\text{wf}}(\tilde{\mu})$, $\tilde{T} = \tilde{T}_{\text{rc}}(\tilde{\mu})$.

Step I. Denote, for $c = 0.11 \ln 2$,

$$\mathcal{A}_{\text{bd}} = \left\{ \tilde{\mu} \in \mathcal{A} : \tilde{t}_{\text{wf}}(\tilde{\mu}) \in [1, c^{-1}], \tilde{T}_{\text{rc}}(\tilde{\mu}) \in [c, 1], \text{ and } \tilde{t}_{\text{wf}}(\tilde{\mu}) \cdot \tilde{T}_{\text{rc}}(\tilde{\mu}) \leq 1 \right\}.$$

As a first step, we show

$$\sup_{\tilde{\mu} \in \mathcal{A}} \Phi(\tilde{\mu}) \leq \max \left\{ 0.11, \sup_{\tilde{\mu} \in \mathcal{A}_{\text{bd}}} \Phi(\tilde{\mu}) \right\}.$$

Let $\kappa = \kappa(\tilde{\mu}) \triangleq \int_{\tilde{t}}^{+\infty} d\tilde{\mu}(x)$ be the mass $\tilde{\mu}$ puts to the waterfilling-active modes.⁶ Expanding $\int \min\{x, \tilde{t}\} d\tilde{\mu}(x) = \int_0^{\tilde{t}} x d\tilde{\mu}(x) + \kappa \tilde{t}$, from Eq. (Dist) we conclude $\tilde{t} \geq 1$, $\tilde{T} \leq 1$, and $\tilde{t}\tilde{T} \leq 1$, where the last inequality is derived as follows:

$$1 = \int \frac{x}{1+x\tilde{T}} d\tilde{\mu}(x) \leq \int_0^{\tilde{t}} x d\tilde{\mu}(x) + \frac{\kappa}{\tilde{T}} = (1 - \kappa\tilde{t}) + \frac{\kappa}{\tilde{T}} = 1 + \kappa \left(-\tilde{t} + \frac{1}{\tilde{T}} \right).$$

We now obtain bounds on $\Phi(\tilde{\mu})$ in terms of \tilde{t}, \tilde{T} :

$$\Phi(\tilde{\mu}) = \underbrace{\frac{1}{2} \int_0^{\tilde{t}} \log(1+x\tilde{T}) d\tilde{\mu}(x)}_{\text{I}} + \underbrace{\frac{1}{2} \int_{\tilde{t}}^{+\infty} \log\left(\frac{1+x\tilde{T}}{x/\tilde{t}}\right) d\tilde{\mu}(x)}_{\text{II}}.$$

From $\int \min\{x, \tilde{t}\} d\tilde{\mu}(x) = 1$ and $\log(1+x) \leq x/\ln 2$, we bound I $\leq \frac{(1-\kappa)\tilde{T}}{2\ln 2}$. To bound II, notice that for $x \geq \tilde{t}$, $\log\left(\frac{1+x\tilde{T}}{x/\tilde{t}}\right) = \log(\tilde{t}\tilde{T} + \tilde{t}/x) \leq \log(\tilde{t}\tilde{T} + 1) \leq \tilde{t}\tilde{T}/\ln 2$, so overall,

$$\text{I} + \text{II} \leq \frac{(1-\kappa)\tilde{T}}{2\ln 2} + \frac{\kappa\tilde{t}\tilde{T}}{2\ln 2} \leq \frac{\tilde{T}}{2\ln 2} + \frac{\tilde{T}}{2\ln 2} = \frac{\tilde{T}}{\ln 2} \leq \frac{1}{\tilde{t}\ln 2}, \quad (34)$$

⁶ $\kappa > 0$ from the water-filling distortion condition.

where the last inequality used previously derived $\tilde{t}\tilde{T} \leq 1$. Then, if $\tilde{t} \geq c^{-1}$ or $\tilde{T} \leq c$, we obtain the desired bound $\Phi(\tilde{\mu}) \leq 0.11$, which concludes this step.

Step II. Second, we show that the largest gap is obtained at spectra, whose mass in the active tail (i.e., values $\geq \tilde{t}$) is distributed between endpoints \tilde{t} and $+\infty$. Concretely, we show that

$$\sup_{\tilde{\mu} \in \mathcal{A}_{\text{bd}}} \Phi(\tilde{\mu}) \leq \sup_{\tilde{\mu} \in \mathcal{A}_{\text{tail}}} \Phi(\tilde{\mu}),$$

where

$$\mathcal{A}_{\text{tail}} = \left\{ \tilde{\mu} \in \mathcal{A}_{\text{bd}} : \text{for } x \geq \tilde{t}_{\text{wf}}(\tilde{\mu}), \theta \in [0, 1], \tilde{\mu}(x) \propto \theta \cdot \delta_{\tilde{t}}(x) + (1 - \theta) \cdot \delta_{+\infty}(x) \right\}.$$

Recall the contribution to the rate gap of the active modes (those $\geq \tilde{t}$) is

$$\frac{1}{2} \int_{\tilde{t}}^{+\infty} \log \left(\frac{1 + x\tilde{T}}{x/\tilde{t}} \right) d\tilde{\mu}(x) = \frac{1}{2} \int_{\tilde{t}}^{+\infty} \log \left(\tilde{t} \cdot \left(\frac{x}{1 + x\tilde{T}} \right)^{-1} \right) d\tilde{\mu}(x).$$

The distortion constraints in Eq. (Dist) yield an equality $\int_{\tilde{t}}^{\infty} \frac{x}{1+x\tilde{T}} d\tilde{\mu}(x) = 1 - \int_0^{\tilde{t}} \frac{x}{1+x\tilde{T}} d\tilde{\mu}(x)$. Note that for $x \in [\tilde{t}, +\infty]$, $\frac{x}{1+x\tilde{T}} \in \left[\frac{\tilde{t}}{1+\tilde{t}\tilde{T}}, \frac{1}{\tilde{T}} \right]$ and is a monotone function. Moreover, both the objective $f(y) = \log(\tilde{t}y^{-1})$ and the first moment constraint are convex, so the supremum is obtained at

$$\tilde{\mu}(x) \propto \theta \cdot \delta_{\tilde{t}}(x) + (1 - \theta) \cdot \delta_{+\infty}(x) \quad \text{for } x \geq \tilde{t}$$

for some $\theta \in [0, 1]$. Note that this suggests that the extremal regime occurs at vanishing distortions.

Step III. Finally, we show that the inactive modes (those $\leq \tilde{t}$) of the worst-case spectra $\tilde{\mu}$ equalize, i.e.,

$$\sup_{\tilde{\mu} \in \mathcal{A}_{\text{tail}}} \Phi(\tilde{\mu}) \leq \sup_{\tilde{\mu} \in \mathcal{A}_{2\text{pt}}} \Phi(\tilde{\mu}),$$

where

$$\mathcal{A}_{2\text{pt}} = \left\{ \tilde{\mu} \in \mathcal{A}_{\text{tail}} : \text{for } \theta \in [0, 1], \tilde{\mu}_0 \leq \tilde{t}_{\text{wf}}(\tilde{\mu}), \tilde{\mu}(x) \propto \theta \cdot \delta_{\tilde{\mu}_0}(x) + (1 - \theta) \cdot \delta_{+\infty}(x) \right\}.$$

For a given $\tilde{\mu} \in \mathcal{A}_{\text{tail}}$ with $\tilde{\mu} \propto (1 - \theta) \cdot \delta_{+\infty}(x)$ for $x > \tilde{t}$, define a corresponding measure $\nu \in \mathcal{A}_{2\text{pt}}$ as

$$\nu \propto \theta \cdot \delta_{\tilde{\mu}_0} + (1 - \theta) \cdot \delta_{+\infty}(x), \quad \text{where } \tilde{\mu}_0 = \frac{\int_0^{\tilde{t}} x d\tilde{\mu}(x)}{\theta}.$$

To complete Step III, it is sufficient to show that $\Phi(\tilde{\mu}) \leq \Phi(\nu)$. Define a path $\mu_s = s \cdot \tilde{\mu} + (1 - s) \cdot \nu$ for $s \in [0, 1]$, and $J(s) \triangleq \Phi(\mu_s)$. We will show that $\frac{\partial J(s)}{\partial s} \leq 0$ for all $s \in [0, 1]$, yielding $\Phi(\tilde{\mu}) = J(1) \leq J(0) = \Phi(\nu)$.

For this, we compute:

$$\frac{\partial \tilde{T}_{\text{rc}}(\mu_s)}{\partial s} = - \frac{\frac{\partial \int \frac{x}{1+x\tilde{T}} d\mu_s(x)}{\partial s}}{\frac{\partial \int \frac{x}{1+x\tilde{T}} d\mu_s(x)}{\partial \tilde{T}}} = \frac{\int \frac{x}{1+x\tilde{T}} d(\tilde{\mu}(x) - \nu(x))}{\int \frac{x^2}{(1+x\tilde{T})^2} d\mu_s(x)}.$$

Note that for all $s \in [0, 1]$, $\tilde{t}_{\text{wf}}(\mu_s)$ is unchanged. Then,

$$\begin{aligned} \frac{\partial 2J(s)}{\partial s} &= \frac{\partial \left(\int_0^{\tilde{t}} \log(1 + x\tilde{T}) d\mu_s(x) + (1 - \theta) \log(\tilde{T}\tilde{t}) \right)}{\partial s} \\ &= \int_0^{\tilde{t}} \log(1 + x\tilde{T}) d(\tilde{\mu}(x) - \nu(x)) + \frac{1}{\ln 2} \cdot \left(\int_0^{\tilde{t}} \frac{x}{1 + x\tilde{T}} d\mu_s(x) + \frac{1 - \theta}{\tilde{T}} \right) \cdot \frac{\int \frac{x}{1 + x\tilde{T}} d(\tilde{\mu}(x) - \nu(x))}{\int \frac{x^2}{(1 + x\tilde{T})^2} d\mu_s(x)} \\ &= \int_0^{\tilde{t}} \left(\log(1 + x\tilde{T}) + A_s \cdot \frac{x}{1 + x\tilde{T}} \right) d(\tilde{\mu}(x) - \nu(x)), \end{aligned}$$

where $A_s \geq 0$. The function $\log(1 + x\tilde{T}) + A_s \cdot \frac{x}{1 + x\tilde{T}}$ is concave on $x \in [0, \tilde{t}]$, yielding the final inequality $\frac{\partial J(s)}{\partial s} \leq 0$. It remains to show $\sup_{\tilde{\mu} \in \mathcal{A}_{2\text{pt}}} \Phi(\tilde{\mu}) \leq 0.11$. Recall that every $\tilde{\mu} \in \mathcal{A}_{2\text{pt}}$ can be expressed for some $\tilde{\mu}_0 \leq \tilde{t}_{\text{wf}}(\tilde{\mu})$ as

$$\tilde{\mu}(x) = \theta \cdot \delta_{\tilde{\mu}_0}(x) + (1 - \theta) \cdot \delta_{+\infty}(x).$$

We can then express

$$2\Phi(\tilde{\mu}) = \theta \log(1 + \tilde{\mu}_0\tilde{T}) + (1 - \theta) \log(\tilde{t}\tilde{T}),$$

where from Eq. (Dist),

$$1 = \theta\tilde{\mu}_0 + (1 - \theta)\tilde{t} = \frac{\theta\tilde{\mu}_0}{1 + \tilde{\mu}_0\tilde{T}} + \frac{1 - \theta}{\tilde{T}}.$$

The case $\theta = 1$ is impossible for $\tilde{T} > 0$, so we will focus on the $\theta \in [0, 1)$ case. We now further simplify this optimization problem to obtain the final 0.11 bound.

From the first distortion constraint, we express $\tilde{t} = (1 - \theta)^{-1} - ((1 - \theta)^{-1} - 1)\tilde{\mu}_0$. From the second distortion constraint, $\tilde{T} = \frac{\theta\tilde{\mu}_0\tilde{T}}{1 + \tilde{\mu}_0\tilde{T}} + 1 - \theta$. Denoting $\alpha = \tilde{\mu}_0\tilde{T}$, we then simplify $\tilde{t}\tilde{T} = (1 - \theta)^{-1}\tilde{T} - ((1 - \theta)^{-1} - 1)\alpha$. Plugging into the rate-gap expression:

$$\begin{aligned} 2\Phi(\tilde{\mu}) &\leq \theta \log(1 + \alpha) + (1 - \theta) \log \left(1 + \frac{\theta\alpha}{(1 + \alpha)(1 - \theta)} - \frac{\alpha\theta}{1 - \theta} \right) \\ &= \theta \log(1 + \alpha) + (1 - \theta) \log \left(1 - \frac{\alpha^2\theta}{(1 + \alpha)(1 - \theta)} \right), \end{aligned}$$

where from the positivity of \tilde{t}, \tilde{T} we have $0 \leq \theta < \frac{1 + \alpha}{1 + \alpha + \alpha^2}$ and $0 \leq \alpha = \tilde{\mu}_0\tilde{T} \leq \tilde{t}\tilde{T} \leq 1$. Denoting $\rho = \frac{\theta}{1 - \theta}$ and differentiating with respect to α gives

$$\begin{aligned} \ln 2(1 - \theta)^{-1} \frac{\partial(2\Phi(\tilde{\mu}))}{\partial \alpha} &= \frac{\rho}{1 + \alpha} - \frac{1}{1 - \frac{\alpha^2\rho}{1 + \alpha}} \cdot \rho \cdot \frac{2\alpha + \alpha^2}{(1 + \alpha)^2} \\ &= \frac{\rho}{1 + \alpha} \left(1 - \frac{2\alpha + \alpha^2}{1 + \alpha - \alpha^2\rho} \right). \end{aligned}$$

Setting this to 0 yields $\rho = \frac{1 - \alpha - \alpha^2}{\alpha^2}$. At the boundary points $\alpha \in \left\{ 0, \frac{\sqrt{5} - 1}{2} \right\}$, the objective function is below the objective of the Theorem. Plugging into the objective, we derive a one-parameter

function

$$2\Delta^* \leq \frac{1 - \alpha - \alpha^2}{1 - \alpha} \log(1 + \alpha) + \frac{\alpha^2}{1 - \alpha} \log\left(1 - \frac{1 - \alpha - \alpha^2}{1 + \alpha}\right).$$

The maximum of this objective is achieved at $\alpha \approx 0.35$ and is $\approx 2 \cdot 0.108 < 2 \cdot 0.11$, which concludes the proof. \blacksquare

Appendix E. Proof of Lemma 14

Lemma 14 (τ Perturbation) *Let R^* be a fixed constant, $C_W = C_W(n) \in \mathbb{R}$, $W \in \mathbb{R}^n$, $s \in \mathbb{R}_+^n$, such that $\|W\|_\infty \leq C_W$ and $\sum_j s_j = n$. Denote*

$$\tau(W, s) = \left(T \sum_j \frac{W_j^2 s_j^2}{(1 + s_j T)^2} \right)^{1/2} \left(\sum_j \frac{s_j}{1 + s_j T} \right)^{-1/2},$$

where $T = T(s)$ is a unique solution to $R_{\text{rc}}^s(T) = R^*$ (see Eq. (RDRC) for definition of $R_{\text{rc}}^s(T)$). If $W, \widehat{W}, s, \widehat{s} \in \mathbb{R}^n$ satisfy $\|W\|_\infty, \|\widehat{W}\|_\infty \leq C_W$ and $\sum_j s_j = \sum_j \widehat{s}_j = n$ and

$$\|W - \widehat{W}\|_\infty \leq \gamma_1 C_W \quad \text{and} \quad \|s - \widehat{s}\|_2 \leq \gamma_2 \sqrt{n},$$

then

$$|\tau(W, s) - \tau(\widehat{W}, \widehat{s})| \leq \gamma_1 C_W + c \cdot C_W \gamma_2 n^4 2^{4nR^*}$$

for sufficiently large n and a constant $c = c(R^*)$.

Proof of Lemma 14. Denote T, \widehat{T} to be the solutions to $R_{\text{rc}}^s(T) = R^*$ and $R_{\text{rc}}^{\widehat{s}}(\widehat{T}) = R^*$ respectively, where recall that $R_{\text{rc}}^{V, s} = R_{\text{rc}}^s$ does not depend on V .

Step 1: Bound on $|T - \widehat{T}|$. First, we show that under the Lemma conditions,

$$|T - \widehat{T}| \leq \gamma_2 n \cdot 2^{4nR^*}. \quad (35)$$

First, since for each $j \in [n]$, $|\log(1 + s_j T) - \log(1 + \widehat{s}_j T)| \leq \frac{1}{\ln 2} \cdot \frac{T}{1 + \min(s_j, \widehat{s}_j) T} \cdot |s_j - \widehat{s}_j| \leq T |s_j - \widehat{s}_j| / \ln 2$, we have

$$\left| \underbrace{R_{\text{rc}}^{\widehat{s}}(\widehat{T})}_{R^*} - R_{\text{rc}}^{\widehat{s}}(T) \right| = |R_{\text{rc}}^s(T) - R_{\text{rc}}^{\widehat{s}}(T)| \leq \frac{T \|s - \widehat{s}\|_1}{2n \ln 2} \leq \frac{T \|s - \widehat{s}\|_2}{2\sqrt{n} \ln 2} \leq \frac{T \gamma_2}{2 \ln 2}.$$

By the mean value theorem, we have for some $\bar{T} \in [\min(T, \widehat{T}), \max(T, \widehat{T})]$,

$$|T - \widehat{T}| = \frac{|R_{\text{rc}}^{\widehat{s}}(\widehat{T}) - R_{\text{rc}}^{\widehat{s}}(T)|}{(R_{\text{rc}}^{\widehat{s}})'(\bar{T})} \leq \frac{T \gamma_2}{2 \ln 2 \cdot (R_{\text{rc}}^{\widehat{s}})'(\bar{T})}.$$

A direct calculation gives $(R_{\text{rc}}^{\widehat{s}})'(\bar{T}) = \frac{1}{2n \ln 2} \sum_j \frac{\widehat{s}_j}{1 + \widehat{s}_j \bar{T}} \geq \frac{1}{2n \ln 2 (1 + \bar{T})} \geq \frac{1}{2n \ln 2 (1 + \max(T, \widehat{T}))}$, where we used that $\max_j \widehat{s}_j \geq 1$. Moreover, for any $s \in \mathbb{R}_+^n$ with $\|s\|_1 = n$, $2n R_{\text{rc}}^s(T) \geq \log(1 + T)$, and therefore, $1 + \max(T, \widehat{T}) \leq 2^{2nR^*}$. This yields

$$|T - \widehat{T}| \leq \gamma_2 n \cdot T (1 + \max(T, \widehat{T})) \leq \gamma_2 n \cdot 2^{4nR^*}.$$

Step 2: Bound on $|\tau(W, s) - \tau(\widehat{W}, s)|$. Here we show that

$$|\tau(W, s) - \tau(\widehat{W}, s)| \leq \gamma_1 C_W. \quad (36)$$

By triangle inequality and $\frac{s_j}{1+s_j T} \leq \frac{1}{T}$,

$$\begin{aligned} |\tau(W, s) - \tau(\widehat{W}, s)| &\leq \left(T \sum_j \frac{(W_j - \widehat{W}_j)^2 s_j^2}{(1 + s_j T)^2} \right)^{1/2} \left(\sum_j \frac{s_j}{1 + s_j T} \right)^{-1/2} \\ &\leq \gamma_1 C_W \cdot \left(T \sum_j \frac{s_j^2}{(1 + s_j T)^2} \right)^{1/2} \left(\sum_j \frac{s_j}{1 + s_j T} \right)^{-1/2} \leq \gamma_1 C_W. \end{aligned}$$

Step 3: Bound on $|\tau(W, s) - \tau(W, \widehat{s})|$. Here we bound $|\tau(W, s) - \tau(W, \widehat{s})|$. For convenience, denote $d_j = \frac{s_j}{1+s_j T} A(s, T) =: \sum_j W_j^2 d_j^2$, $B(s, T) =: \sum_j d_j$ and let $A = A(s, T)$, $\widehat{A} = A(\widehat{s}, \widehat{T})$ and $B = B(s, T)$, $\widehat{B} = B(\widehat{s}, \widehat{T})$, so that

$$|\tau(W, s) - \tau(W, \widehat{s})| = \left| \frac{\sqrt{TA}}{\sqrt{B}} - \frac{\sqrt{\widehat{T}\widehat{A}}}{\sqrt{\widehat{B}}} \right| \leq \sqrt{T} \left| \frac{\sqrt{A}}{\sqrt{B}} - \frac{\sqrt{\widehat{A}}}{\sqrt{\widehat{B}}} \right| + |\sqrt{T} - \sqrt{\widehat{T}}| \cdot \frac{\sqrt{\widehat{A}}}{\sqrt{\widehat{B}}} =: \text{I} + \text{II}$$

by triangle inequality. To bound I, we use a triangle inequality $\text{I} = \sqrt{T} \left| \frac{\sqrt{A}}{\sqrt{B}} - \frac{\sqrt{\widehat{A}}}{\sqrt{\widehat{B}}} \right| \leq \sqrt{T} \frac{|\sqrt{A} - \sqrt{\widehat{A}}|}{\sqrt{B}} + \sqrt{T\widehat{A}} \left| \frac{1}{\sqrt{B}} - \frac{1}{\sqrt{\widehat{B}}} \right|$ and notice

$$|\sqrt{A} - \sqrt{\widehat{A}}| = \left| \|W \odot d\|_2 - \|W \odot \widehat{d}\|_2 \right| \leq \|W\|_\infty \|d - \widehat{d}\|_2 \leq C_W \|d - \widehat{d}\|_2.$$

Moreover,

$$|\sqrt{B} - \sqrt{\widehat{B}}| = \left| \|d\|_1^{1/2} - \|\widehat{d}\|_1^{1/2} \right| \leq \frac{\sqrt{n} \|d - \widehat{d}\|_2}{\sqrt{B} + \sqrt{\widehat{B}}}.$$

For a uniform lower bound on B, \widehat{B} we use $\min(B, \widehat{B}) \geq \frac{1}{1+T}$, so the combined bound is

$$\begin{aligned} \text{I} &\leq \sqrt{T} \|d - \widehat{d}\|_2 \left[C_W (1+T)^{1/2} + \sqrt{\widehat{A}n} (1+T)^{3/2} \right] \\ &\leq 2 \|d - \widehat{d}\|_2 C_W n 2^{2nR^*}. \end{aligned}$$

where we used $\widehat{A} \leq C_W^2 \frac{n}{T^2}$ and $1+T \leq 2^{2nR^*}$. It remains to bound $\|d - \widehat{d}\|_2$. Denoting $d(s, T) = \left(\frac{s_j}{1+s_j T} \right)_j$ (so that $d = d(s, T)$ and $\widehat{d} = d(\widehat{s}, \widehat{T})$), we have

$$\|d - \widehat{d}\|_2 \leq \|d(s, T) - d(\widehat{s}, T)\|_2 + \|d(\widehat{s}, T) - d(\widehat{s}, \widehat{T})\|_2.$$

Since $\left| \frac{s_j}{1+s_j T} - \frac{\widehat{s}_j}{1+\widehat{s}_j T} \right| \leq |s_j - \widehat{s}_j|$, we have $\|d(s, T) - d(\widehat{s}, T)\|_2 \leq \|d(s, T) - d(\widehat{s}, T)\|_1 \leq \|s - \widehat{s}\|_1 \leq \gamma_2 n$. Moreover, $\left| \frac{s_j}{1+s_j T} - \frac{s_j}{1+s_j \widehat{T}} \right| \leq n^2 |T - \widehat{T}|$, so we obtain

$$\|d - \widehat{d}\|_2 \leq \gamma_2 n + 2\gamma_2 n^4 2^{4nR^*} \leq 3\gamma_2 n^4 2^{4nR^*},$$

where we used $|T - \widehat{T}| \leq 2\gamma_2 n 2^{4nR^*}$ in Eq. (35). To bound II, we use the bound in Eq. (35) and the lower bound $T \geq 2R^* \ln 2$ obtained in the proof of Thm. 7:

$$|\sqrt{T} - \sqrt{\widehat{T}}| = \frac{|T - \widehat{T}|}{\sqrt{T} + \sqrt{\widehat{T}}} \leq \frac{\gamma_2 n 2^{4nR^*}}{\sqrt{2R^* \ln 2}}.$$

Then,

$$\text{II} \leq \frac{\gamma_2 n 2^{4nR^*}}{\sqrt{2R^* \ln 2}} \cdot \frac{C_W \sqrt{n}}{2T} \cdot \sqrt{1+T} \leq \gamma_2 C_W n^{3/2} 2^{4nR^*} / (2R^* \ln 2).$$

We have for the final bound

$$\begin{aligned} |\tau(W, s) - \tau(W, \widehat{s})| &\leq \text{I} + \text{II} \leq 6\gamma_2 n^{4/2} 2^{4nR^*} C_W n 2^{2nR^*} + \gamma_2 C_W n^{3/2} 2^{4nR^*} / (2R^*) \\ &\leq c \cdot C_W \gamma_2 n^{4/2} 2^{4nR^*} \end{aligned}$$

for sufficiently large n and a constant $c = c(R^*)$. Analogous bound holds for $|\tau(\widehat{W}, s) - \tau(\widehat{W}, \widehat{s})|$.

Step 4: Final bound. Combining Steps 2,3 above, we obtain

$$\begin{aligned} |\tau(W, s) - \tau(\widehat{W}, \widehat{s})| &\leq |\tau(W, s) - \tau(\widehat{W}, s)| + |\tau(\widehat{W}, s) - \tau(\widehat{W}, \widehat{s})| \\ &\leq \gamma_1 C_W + c \cdot C_W \gamma_2 n^{4/2} 2^{4nR^*}. \end{aligned}$$

□

Appendix F. Proof of upper bound in oracle Proposition 2

Here we sketch proof of (3). First, since f and g are allowed to depend on Σ_X we can rotate Σ_X to eigenbasis, and thus assume from now on that $\Sigma_X = \text{diag}(\lambda_1, \dots, \lambda_n)$. Fix t and let $D_i = \min\{t/\lambda_i, 1\}$. Let $R_0 = R_{\text{wf}}(\Sigma_X, t)$. Fix arbitrary $\epsilon > 0$ and set rate $R = R_0 + 2\epsilon$. We will show that by generating codebook \mathbf{C} randomly via sampling $1 + 2^{nR}$ codewords from distribution $\prod_{i=1}^n \mathcal{N}(0, 1 - D_i)$ one can attain distortion

$$\mathbb{E}_W \left[\min_{c \in \mathbf{C}} d_{\Sigma_X}(W, c) \right] \leq n D_{\text{wf}}(\Sigma_X, t) + \left(e^{-2n\epsilon} + c_1 e^{-c_2 n \epsilon^2} \right) \text{tr} \Sigma_X, \quad (37)$$

where $c_1, c_2 > 0$ are some absolute constants. From here the statement of the theorem follows by setting $\epsilon_n = c \sqrt{\frac{\log n}{n}}$ with an appropriate $c > 0$.

To show (37) let $Y_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1 - D_i)$ independently. Let also $Z_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$. Set

$$W_i = Y_i + \sqrt{D_i} Z_i$$

and notice that $W = (W_1, \dots, W_n) \sim \mathcal{N}(0, I_n)$. This coupling of W to Y satisfies a useful property:

$$\mathbb{E}_W [d_{\Sigma_X}(W, Y)] = \mathbb{E}_W \left[\sum_{i=1}^n \lambda_i D_i Z_i^2 \right] = D_{\text{wf}}(\Sigma_X, t).$$

On the other hand, we have *information density*

$$i(W; Y) := \log \frac{dP_{W,Y}}{d(P_W \times P_Y)}(W, Y) = nR_0 + \frac{\log e}{2} \sum_{i=1}^n W_i^2 - Z_i^2.$$

Note that $\mathbb{E}[i(W; Y)] = I(W; Y) = nR_0$. Note that while W_i, Z_i are jointly correlated Gaussians, they are independent for different i 's. A fascinating property of information density (underlying information stability) of Gaussian processes is that its variance is uniformly bounded regardless of the distribution, see (Polyanskiy and Wu, 2024, (19.32)). We will exploit this below to show “local” subgaussian estimate on the concentration of $i(W; Y)$.

Define log-MGF function

$$f(z) := \frac{1}{n} \ln \mathbb{E}_{W,Z} \left[e^{z \sum_{i=1}^n W_i^2 - Z_i^2} \right] = \frac{1}{n} \sum_{i=1}^n \ln \mathbb{E}_{W_i, Z_i} \left[e^{z(W_i^2 - Z_i^2)} \right].$$

It is not hard to show that there exists a neighborhood \mathcal{S} of 0 on the complex plane \mathbb{C} and a constant c such that

$$\sup_{z \in \mathcal{S}} |f(z)| \leq c,$$

and crucially \mathcal{S} and c can be chosen independent of $\{D_i\}$ (and hence of t, λ_i). For this, one only needs to apply Cauchy-Schwarz to reduce to analysis of log-MGF of W_i^2 and Z_i^2 , which are just squares of $\mathcal{N}(0, 1)$.

In addition to being analytic, f also satisfies $f(0) = f'(0) = 0$ and $f''(0) \leq 2\mathbb{E}[W_i^4 + Z_i^4] \leq 12$. Thus, by Cauchy formula we can uniformly bound $f''(z)$ inside any compact subset of \mathcal{S} . Consequently, there exists a universal $c' > 0$ and $z_0 > 0$ such that for all real $-z_0 < z < z_0$ we have

$$f(z) \leq 2c' z^2.$$

Applying Chernoff estimate we find that for all $\epsilon < \epsilon_0$ we have

$$\mathbb{P}[i(W; Y) > I(W; Y) + n\epsilon] \leq e^{-2c_2 n \epsilon^2}, \quad (38)$$

where crucially ϵ_0, c_2 are absolute constants.

We are now ready to apply standard finite blocklength rate-distortion upper bound (Polyanskiy and Wu, 2024, Theorem 25.2), which claims existence of codebook \mathbf{C} with distortion

$$\mathbb{E}_W \left[\min_{c \in \mathbf{C}} d_{\Sigma_X}(W, c) \right] \leq \mathbb{E}[d_{\Sigma_X}(W, Y)] + \mathbb{E}[d_{\Sigma_X}(W, 0)] e^{-2nR/\gamma} + \mathbb{E}[d_{\Sigma_X}(W, Y) 1\{i(W; Y) > \log_2 \gamma\}], \quad (39)$$

where γ is arbitrary, but we set it to

$$\log_2 \gamma = nR_0 + n\epsilon.$$

Applying Cauchy-Schwarz and (38) to the last term in (39) we obtain (37).

Appendix G. Additive rate-distortion for quantization of a colored X

It is worth mentioning that the rate-distortion region (RDRC) we obtained also characterizes the tradeoff between rate and distortion of a particular quantization scheme in a different, but closely related setup.

In particular, let $X \sim \mathcal{N}(0, \Sigma_X)$ be a Gaussian vector in \mathbb{R}^n to be quantized under the standard quadratic distortion measure $D = \frac{1}{n} \mathbb{E} \|\hat{X} - X\|_2^2$. Clearly the optimal rate-distortion tradeoff for this problem is given by

$$R(D) = \frac{1}{n} \min I(X; \hat{X}) \quad (40)$$

where the minimum is over all $P_{\hat{X}|X}$ for which $\frac{1}{n} \mathbb{E} \|\hat{X} - X\|_2^2 \leq D$. As we already discussed, the optimal $P_{\hat{X}|X}$ is determined by the reverse waterfilling solution, and is given precisely by (WF). In fact, we derived the oracle lower bound by showing that in the oracle setup, where the decoder also knows Σ_X , our problem is equivalent to that of quantizing X under standard quadratic loss.

While the waterfilling solution gives the optimal tradeoff $R^*(D)$ function, any other valid choice of $P_{\hat{X}|X}$ yields an achievable $R(D)$. A very simple choice is to construct \hat{X} by first adding independent noise $Z \sim \mathcal{N}(0, \frac{1}{T} I_n)$ to X and then performing minimum mean squared error (MMSE) estimation of X from $X + Z$. The $R(D)$ tradeoff attained by this particular choice of $P_{\hat{X}|X}$ has been studied in the information theory literature (for sources that are not necessarily Gaussian) under the name *additive rate-distortion function (ARDF)* Zamir (2002); Zamir and Berger (2002); Ostergaard and Zamir (2011).

The MMSE estimator of X from $X + Z$ is linear and therefore $\hat{X} = F(X + Z)$, where

$$F = \Sigma_X \left(\Sigma_X + \frac{1}{T} I_n \right)^{-1} = U \cdot \text{diag} \left(\frac{\lambda_1 T}{1 + \lambda_1 T}, \dots, \frac{\lambda_n T}{1 + \lambda_n T} \right) \cdot U^\top \quad (41)$$

and

$$\frac{1}{n} \mathbb{E} \|\hat{X} - X\|_2^2 = \frac{1}{n} \sum_{i=1}^n \frac{\lambda_i}{1 + \lambda_i T} = D_{\text{rc}}(\lambda, T). \quad (42)$$

Furthermore, if all singular values of Σ_X are positive, F is invertible and

$$\frac{1}{n} I(X; \hat{X}) = \frac{1}{n} I(X; X + Z) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \log(1 + \lambda_i T) = R_{\text{rc}}(\lambda, T). \quad (43)$$

Thus, the Σ_X universal rate distortion tradeoff we derived for the problem of quantizing a white source $\mathcal{N}(0, I_n)$ under d_{Σ_X} metric known only to the encoder is precisely the additive rate-distortion function for quantizing $X \sim \mathcal{N}(0, \Sigma_X)$ under standard quadratic loss.