

# Mutual Information as a Function of Matrix SNR for Linear Gaussian Channels

Galen Reeves<sup>\*†</sup>, Henry D. Pfister<sup>\*</sup> and Alex Dytso<sup>‡</sup>

<sup>\*</sup>Department of Electrical Engineering, Duke University

<sup>†</sup>Department of Statistical Science, Duke University

<sup>‡</sup>Department of Electrical Engineering, Princeton University

**Abstract**—This paper focuses on the mutual information and minimum mean-squared error (MMSE) as a function of a matrix-valued signal-to-noise ratio (SNR) for a linear Gaussian channel with arbitrary input distribution. As shown by Lamarca, the mutual-information is a concave function of a positive semi-definite matrix, which we call the matrix SNR. This implies that the mapping from the matrix SNR to the MMSE matrix is decreasing monotone. Building upon these functional properties, we start to construct a unifying framework that provides a bridge between classical information-theoretic inequalities, such as the entropy power inequality, and interpolation techniques used in statistical physics and random matrix theory. This framework provides new insight into the structure of phase transitions in coding theory and compressed sensing. In particular, it is shown that the parallel combination of linear channels with freely-independent matrices can be characterized succinctly via free convolution.

**Index Terms**—I-MMSE, entropy power inequality, conditional central limit theorem, random matrix theory, compressed sensing, Gaussian logarithmic Sobolev inequality.

## I. INTRODUCTION

The functional properties of mutual information play an important role in applications across the mathematical sciences. In some cases, these functional properties lead to simple and elegant proofs of deep mathematical results. In other cases, they provide a crucial step in the proof of new results.

The focus of this paper is on the mutual information and minimum mean-squared error (MMSE) associated with the linear Gaussian channel, which is described by

$$\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{N}, \quad (1)$$

where  $\mathbf{A}$  is a  $k \times n$  channel matrix,  $\mathbf{X}$  is an  $n$ -dimensional random vector (or signal), and  $\mathbf{N} \sim \mathcal{N}(0, I_n)$  is standard Gaussian noise. Our motivation comes largely from the fact that  $I(\mathbf{X}; \mathbf{A}\mathbf{X} + \mathbf{N})$  is closely related to interesting open questions in compressed sensing, coding theory, and statistical physics. In particular, the inference problems associated with these applications can all have phase transitions as the dimension increases.

For example, if  $\mathbf{X}$  is uniformly distributed on points in a codebook and  $\mathbf{A} = \sqrt{s}\mathbf{I}$ , then we have coded communication

The work of G. Reeves was supported in part by funding from the Laboratory for Analytic Sciences (LAS). The work of G. Reeves and H. Pfister was supported part by the NSF under Grant No. 1718494. Any opinions, findings, conclusions, and recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the sponsors.

over a Gaussian channel and the error probability can have a phase transition (as the problem size grows) when  $s$  increases [1]. If instead, one chooses  $\mathbf{A}$  to be a random diagonal matrix where each entry is 0 with probability  $\epsilon$  and  $\sqrt{s}$  with probability  $1 - \epsilon$ , then this approximates (for large  $s$ ) an erasure channel and the error probability can have a phase transition as  $\epsilon$  increases [2]. When  $\mathbf{X}$  is a sparse random vector (e.g., i.i.d. with a probability mass at 0), there can be a phase transition in the MMSE as the sparsity level increases [3]. The locations of these phase transitions can sometimes be predicted by the non-rigorous replica method from statistical physics [4]. An interesting open question for all of these problems is determining if and when these replica formulas are correct.

Recently, the first two authors used information-theoretic tools to prove that the replica-symmetric prediction for the MMSE function in compressed sensing is exact when the channel matrix is drawn i.i.d. from the Gaussian ensemble [5]. Subsequent work established a similar result using a different proof technique [6], which was later extended to the generalized linear model [7]. Based on this, the first author described how these ideas might be extended to right-orthogonally invariant channel matrices and multilayer linear models [8]. This paper provides some rigorous progress on this problem.

## A. Matrix SNR

The key insight underlying the results in this paper is that the linear Gaussian channel can be parameterized by a positive semi-definite matrix that generalizes the scalar signal-to-ratio (SNR). This insight follows from the fact that a linear Gaussian channel with  $k \times n$  channel matrix  $\mathbf{A}$  is statistically equivalent to a linear Gaussian channel with  $n \times n$  channel matrix  $(\mathbf{A}^T \mathbf{A})^{\frac{1}{2}}$ , regardless of the signal distribution. Although this fact follows naturally from the independence and orthogonal invariance of the Gaussian noise, our literature review did not reveal any references before the 2009 paper of Lamarca [9], and we believe this result has not been fully exploited. This approach also brings the multivariate derivative formula for the mutual information closer to the well-developed theory associated with the scalar I-MMSE relationship [10].

An important consequence of the parametrization described above is that one obtains an explicit additivity rule for linear Gaussian channels. Specifically, two parallel channels with matrices  $\mathbf{A}$  and  $\mathbf{B}$ , respectively, are equivalent to a channel

with  $n \times n$  matrix  $(A^T A + B^T B)^{\frac{1}{2}}$ . Combining this fact with the chain rule and the low-noise expansion of mutual information [11], one finds that: 1) the mutual information  $I(\mathbf{X}; \mathbf{A}\mathbf{X} + \mathbf{N})$  is concave in  $A^T A$ ; and 2) the minimum mean-square error matrix  $\mathbb{E}[\text{Cov}(\mathbf{X} | \mathbf{A}\mathbf{X} + \mathbf{N})]$  is decreasing monotone map of  $A^T A$ . Mathematically, this means that for any  $n$ -dimensional random vector  $\mathbf{X}$  and matrices  $A, B$  with  $n$  columns, one has

$$\begin{aligned} \lambda I(\mathbf{X}; \mathbf{A}\mathbf{X} + \mathbf{N}) + (1 - \lambda)I(\mathbf{X}; \mathbf{B}\mathbf{X} + \mathbf{N}') \\ \leq I(\mathbf{X}; (\lambda A^T A + (1 - \lambda)B^T B)^{\frac{1}{2}} \mathbf{X} + \mathbf{N}''), \end{aligned} \quad (2)$$

for all  $0 \leq \lambda \leq 1$  where  $\mathbf{N}, \mathbf{N}'$  and  $\mathbf{N}''$  are independent standard Gaussian vectors. Furthermore, the MMSE matrix satisfies  $\langle A^T A - B^T B, \Delta \rangle \leq 0$  where  $\Delta = \mathbb{E}[\text{Cov}(\mathbf{X} | \mathbf{A}\mathbf{X} + \mathbf{N})] - \mathbb{E}[\text{Cov}(\mathbf{X} | \mathbf{B}\mathbf{X} + \mathbf{N}')] and  $\langle \cdot, \cdot \rangle$  denotes the trace inner-product on matrices.$

**Remark 1.** Although it may be tempting to reinterpret (2) in terms of differential entropy, such a decomposition becomes cumbersome due to a mismatch between problem dimensions.

**Remark 2.** It is important to note that (2) holds for arbitrary matrices. Increasing the scale of the matrices  $A \mapsto \sqrt{s}A$  for some positive number  $s$  is equivalent to decreasing the noise power. If  $\mathbf{X}$  has finite entropy then the low-noise limit is well defined and one obtains

$$\begin{aligned} \lambda H(\mathbf{X} | \mathbf{A}\mathbf{X}) + (1 - \lambda)H(\mathbf{X} | \mathbf{B}\mathbf{X}) \\ \geq H(\mathbf{X} | (\lambda A^T A + (1 - \lambda)B^T B)^{\frac{1}{2}} \mathbf{X}). \end{aligned} \quad (3)$$

### B. Our Contributions

Building upon the aforementioned functional properties, we obtain the following results:

- An effective Fisher information matrix is introduced and shown to have a monotonicity property that provides a multivariate version of the single-crossing property (Theorem 2). Using this result, we obtain matrix inequalities that mimic well-known scalar bounds on the mutual information.
- A new relationship is identified between three different measures of the distance between the distribution of an  $n$ -dimensional random vector and an i.i.d. Gaussian distribution. In particular, Theorem 3 shows that the relative entropy measure can be expressed as the sum of relative entropy associated with a (random) low-dimensional projection and a second term that can be related to the deficit in the entropy power inequality (EPI).
- Using ideas from free probability theory, concentration results and functional properties are established that characterize the asymptotic behavior of right-orthogonally invariant random matrices (Theorem 4). Although these results are closely linked to the conjectured limits implied by the replica method, they are proven rigorously here under mild assumptions.

Due to space constraints many of the proofs are either sketched or omitted.

### C. Notation

We use  $\mathbb{S}^n$ ,  $\mathbb{S}_+^n$  and  $\mathbb{S}_{++}^n$  to denote the space  $n \times n$  symmetric matrices, positive semi-definite matrices, positive definite matrices, respectively. Given a positive-definite matrix  $S$ , we use  $S^{\frac{1}{2}}$  to denote the positive-definite square root.

## II. MATRIX-SNR FUNCTIONS

For any  $n$ -dimensional random vector  $\mathbf{X}$ , the mutual information function  $I_{\mathbf{X}} : \mathbb{S}_+^n \rightarrow [0, \infty)$  and MMSE function  $M_{\mathbf{X}} : \mathbb{S}_+^n \rightarrow \mathbb{S}_+^n$  are defined by

$$I_{\mathbf{X}}(S) \triangleq I(\mathbf{X}; S^{\frac{1}{2}} \mathbf{X} + \mathbf{N}) \quad (4)$$

$$M_{\mathbf{X}}(S) \triangleq \mathbb{E}[\text{Cov}(\mathbf{X} | S^{\frac{1}{2}} \mathbf{X} + \mathbf{N})], \quad (5)$$

where  $\mathbf{N}$  is a standard Gaussian vector. The following result shows that these functions can be related to the mutual information and MMSE associated with an arbitrary  $k \times n$  matrix  $A$ . For completeness, we also sketch a partial proof below.

**Lemma 1** ([9]). *For every  $k \times n$  matrix  $A$ , we have*

$$I_{\mathbf{X}}(A^T A) = I(\mathbf{X}; \mathbf{A}\mathbf{X} + \mathbf{W}) \quad (6)$$

$$M_{\mathbf{X}}(A^T A) = \mathbb{E}[\text{Cov}(\mathbf{X} | \mathbf{A}\mathbf{X} + \mathbf{W})], \quad (7)$$

where  $\mathbf{W} \sim \mathcal{N}(0, I_k)$ .

*Sketch of Proof.* First, consider the case  $k < n$ . Let  $B$  be the  $n \times n$  matrix whose first  $k$  rows are equal to  $A$  and whose remaining  $(n - k)$  rows are equal to zero. Clearly,  $I(\mathbf{X}; \mathbf{A}\mathbf{X} + \mathbf{W}) = I(\mathbf{X}; \mathbf{B}\mathbf{X} + \mathbf{N})$ , because the extra rows in  $B$  do not convey any information about  $\mathbf{X}$ . Next, consider the singular value decomposition  $B = U\Sigma V^T$  where all matrices are  $n \times n$ . By the orthogonal invariance of the standard Gaussian distribution, we can write

$$\begin{aligned} I(\mathbf{X}; \mathbf{B}\mathbf{X} + \mathbf{N}) &= I(\mathbf{X}; U\Sigma V^T \mathbf{X} + \mathbf{N}) \\ &= I(\mathbf{X}; \Sigma V^T \mathbf{X} + \mathbf{N}) \\ &= I(\mathbf{X}; V\Sigma V^T \mathbf{X} + \mathbf{N}). \end{aligned}$$

Noting that  $V\Sigma V^T = (B^T B)^{\frac{1}{2}} = (A^T A)^{\frac{1}{2}}$  gives the stated identity. The case of  $k \geq n$  follows from similar arguments and proofs of the MMSE result (7) can be found in [9], [12].  $\square$

**Lemma 2.** *For all  $S, T \in \mathbb{S}_+^n$ , we have*

$$I_{\mathbf{X}}(S + T) = I(\mathbf{X}; \mathbf{Y}, \mathbf{Z}) \quad (8)$$

$$M_{\mathbf{X}}(S + T) = \mathbb{E}[\text{Cov}(\mathbf{X} | \mathbf{Y}, \mathbf{Z})], \quad (9)$$

where  $\mathbf{Y} = S^{\frac{1}{2}} \mathbf{X} + \mathbf{N}$ ,  $\mathbf{Z} = T^{\frac{1}{2}} \mathbf{X} + \mathbf{N}'$ , and  $\mathbf{N}$  and  $\mathbf{N}'$  are independent standard Gaussian vectors.

*Proof.* The right-hand side of (8) can be expressed as  $I(\mathbf{X}; \mathbf{A}\mathbf{X} + \mathbf{N}'')$  where  $A$  is the  $2n \times n$  matrix obtained by stacking  $S^{\frac{1}{2}}$  and  $T^{\frac{1}{2}}$  and  $\mathbf{N}''$  is the  $2n \times 1$  vector obtained by stacking  $\mathbf{N}$  and  $\mathbf{N}'$ . Noting that  $A^T A = S + T$  and invoking Lemma 1 gives the stated result.  $\square$

**Lemma 3** (Scaling property). For any  $k \times n$  matrix  $A$  and  $S \in \mathbb{S}_+^n$ , we have

$$I_{A\mathbf{X}}(S) = I_{\mathbf{X}}(ASA^T) \quad (10)$$

$$M_{A\mathbf{X}}(S) = AM_{\mathbf{X}}(ASA^T)A^T. \quad (11)$$

**Lemma 4.** The mutual information function  $I_{\mathbf{X}}(S)$  is twice differentiable on  $\mathbb{S}_{++}^n$  with gradient and Hessian given by

$$\nabla_S I_{\mathbf{X}}(S) = \frac{1}{2} M_{\mathbf{X}}(S) \quad (12)$$

$$\nabla_S^2 I_{\mathbf{X}}(S) = -\frac{1}{2} \mathbb{E}_{\mathbf{Y}}[\Phi_{\mathbf{X}}(\mathbf{Y}) \otimes \Phi_{\mathbf{X}}(\mathbf{Y})] \quad (13)$$

where  $\mathbf{Y} = S^{\frac{1}{2}}\mathbf{X} + \mathbf{N}$  and  $\Phi_{\mathbf{X}}(\mathbf{y}) = \text{Cov}(\mathbf{X} | \mathbf{Y} = \mathbf{y})$ .

**Remark 3.** The Hessian is negative semi-definite since the covariance  $\Phi_{\mathbf{X}}(\mathbf{Y})$  is positive semi-definite. This implies that (6) is a concave function on  $\mathbb{S}_{++}^n$ . This concavity and the results in Lemma 4 can be found in Lamarca [9] and Payaró et al. [12].

The concavity is summarized by the following theorem.

**Theorem 1.** For any  $n$ -dimensional random vector  $\mathbf{X}$ , The mutual information is concave on  $\mathbb{S}_+^n$ . In other words,

$$(1 - \lambda)I_{\mathbf{X}}(S) + \lambda I_{\mathbf{X}}(T) \geq I_{\mathbf{X}}((1 - \lambda)S + \lambda T) \quad (14)$$

for all  $S, T \in \mathbb{S}_+^n$  and  $\lambda \in [0, 1]$ . Furthermore, the MMSE function (7) is a decreasing monotone mapping on  $\mathbb{S}_+^n$ . Thus,

$$\text{tr}((S - T)(M_{\mathbf{X}}(S) - M_{\mathbf{X}}(T))) \leq 0 \quad (15)$$

for all  $S, T \in \mathbb{S}_+^n$ .

### III. APPLICATIONS

This section describes some new results that can be obtained using properties of the mutual information and MMSE matrix as a function of the matrix SNR.

#### A. Bounds on the mutual information and MMSE

The Fisher information matrix of a random vector  $\mathbf{Y}$  with density  $p(\mathbf{y})$  is defined to be  $J(\mathbf{Y}) \triangleq \text{Cov}(\rho(\mathbf{Y}))$  where  $\rho(\mathbf{y}) \triangleq \nabla \log p(\mathbf{y})$  is the score function. Given any random vector  $\mathbf{X}$  with  $\text{Cov}(\mathbf{X}) \in \mathbb{S}_{++}^n$ , we define the effective Fisher information matrix  $K_{\mathbf{X}}: \mathbb{S}_+^n \rightarrow \mathbb{S}_+^n$  by

$$K_{\mathbf{X}}(S) \triangleq M_{\mathbf{X}}^{-1}(S) - S. \quad (16)$$

This is the Fisher information matrix of the multivariate Gaussian distribution whose MMSE matrix equals  $M_{\mathbf{X}}(S)$ . Using the matrix version of Brown's identity [13, Proposition 6]

$$J(\mathbf{X} + S^{-\frac{1}{2}}\mathbf{N})S^{-1} + SM_{\mathbf{X}}(S) = I, \quad (17)$$

it can be verified that

$$K_{\mathbf{X}}^{-1}(S) = J^{-1}(\mathbf{X} + S^{-\frac{1}{2}}\mathbf{N}) - S^{-1} \quad (18)$$

for all  $S \in \mathbb{S}_{++}^n$ . This characterization shows that  $K_{\mathbf{X}}^{-1}(S)$  is well-defined even if  $\text{Cov}(\mathbf{X})$  is degenerate.

**Lemma 5.** The effective Fisher information matrix  $K_{\mathbf{X}}(S)$  is decreasing monotone on  $\mathbb{S}_+^n$  with

$$K_{\mathbf{X}}(0) = \text{Cov}^{-1}(\mathbf{X}), \quad (19)$$

$$\lim_{\lambda_{\min}(S) \rightarrow \infty} K_{\mathbf{X}}(S) = J(\mathbf{X}). \quad (20)$$

*Sketch of Proof.* As the gradient is negative semi-definite, monotonicity follows. The limits come from (16) and (18).  $\square$

**Theorem 2.** Let  $\mathbf{X}$  be an  $n$ -dimensional random vector with positive-definite covariance matrix. For any matrices  $R, S, T \in \mathbb{S}_+$  with  $R \preceq S \preceq T$ , the MMSE matrix satisfies

$$(K_{\mathbf{X}}(T) + S)^{-1} \preceq M_{\mathbf{X}}(S) \preceq (K_{\mathbf{X}}(R) + S)^{-1}. \quad (21)$$

*Sketch of Proof.* This result follows from the monotonicity of  $K_{\mathbf{X}}(S)$ , established in Lemma 5, and the fact that inversion reverses the partial-order on positive semi-definite matrices.  $\square$

Theorem 2 can be thought of as a matrix generalization of the single-crossing property (see [14]). Taking the limits  $\lambda_{\max}(T) \rightarrow 0$  and  $\lambda_{\min}(R) \rightarrow \infty$ , one obtains the bounds

$$(J(\mathbf{X}) + S)^{-1} \preceq M_{\mathbf{X}}(S) \preceq (\text{Cov}^{-1}(\mathbf{X}) + S)^{-1}. \quad (22)$$

Note that the left-hand side of (22) is a matrix version of the Bayesian Cramer-Rao lower bound [15]. Here the lower bound is meaningful only if  $\mathbf{X}$  has finite Fisher information. The right hand side of (22) is often called the linear MMSE.

By the multivariate I-MMSE relationship (10), we see that the mutual information can be expressed as

$$I_{\mathbf{X}}(S) = \frac{1}{2} \int_0^1 \text{tr} \left( M_{\mathbf{X}}(S_t) \frac{d}{dt} S_t \right) dt \quad (23)$$

where the integral is over any differentiable path  $t \mapsto S_t$  with  $S_0 = 0$  and  $S_1 = S$ . In particular, letting  $S_t = tS$ , swapping the integral and the trace, and using Lemma 3, leads to

$$I_{\mathbf{X}}(S) = \frac{1}{2} \text{tr} \left( S \int_0^1 M_{\mathbf{X}}(tS) dt \right) \quad (24)$$

$$= \frac{1}{2} \text{tr} \left( \int_0^1 M_{S^{\frac{1}{2}}\mathbf{X}}(tI) dt \right). \quad (25)$$

Interestingly, this decomposition shows that the mutual information can be viewed as the trace of an integrated MMSE matrix. A similar observation was made by Dembo [16, pg. 14], who showed that an entropy can be expressed as the trace of an integral involving the Fisher information matrix.

Combining Theorem 2 with (25) provides a lower bound on the *matrix* inside the trace in (25):

$$\begin{aligned} \int_0^1 M_{S^{\frac{1}{2}}\mathbf{X}}(tI) dt &= \int_0^1 \left( K_{S^{\frac{1}{2}}\mathbf{X}}(tI) + tI \right)^{-1} dt \\ &\succeq \int_0^1 \left( K_{S^{\frac{1}{2}}\mathbf{X}}(I) + tI \right)^{-1} dt \\ &= \log \left( I + K_{S^{\frac{1}{2}}\mathbf{X}}^{-1}(I) \right) \\ &= \log \left( J^{-1}(S^{\frac{1}{2}}\mathbf{X} + \mathbf{N}) \right) \\ &\succeq \log \left( I + S^{\frac{1}{2}} J^{-1}(\mathbf{X}) S^{\frac{1}{2}} \right). \end{aligned}$$

Here, we recall that the matrix logarithm is well-defined on  $\mathbb{S}_{++}^n$ . Taking the trace of both sides recovers some well-known lower bounds on the mutual information

$$\begin{aligned} I_{\mathbf{X}}(S) &\geq \frac{1}{2} \log \det(J^{-1}(S^{\frac{1}{2}} \mathbf{X} + \mathbf{N})) \\ &\geq \frac{1}{2} \log \det(I + SJ^{-1}(\mathbf{X})). \end{aligned}$$

### B. Gaussian approximation via low-dimensional projections

Many of the standard approaches to measuring relative entropy with respect to the Gaussian measure are based on the fact that adding an independent Gaussian component to a random vector decreases the relative entropy to a Gaussian (see [17, Section 3.2]).

A related but different phenomenon is that most low-dimensional projections of a high-dimensional distribution are closer to an i.i.d. Gaussian, in a relative sense, than the original distribution. Some recent information-theoretic bounds, known as conditional central limit theorems (CCLTs), are provided by the first author in [18]. One motivation for the present paper is to understand how the decrease in approximation error associated with low-dimensional projections relates to properties of the original distribution.

For concreteness, we will focus on projections on the Stiefel manifold  $\mathcal{V}_k(\mathbb{R}^n)$ , which is the set of  $k \times n$  matrices  $A$  satisfying  $AA^T = I_k$ . Furthermore, we will consider distributions of the form  $\sqrt{s}\mathbf{X} + \mathbf{N}$  where  $s \in (0, \infty)$  and  $\mathbf{N}$  is a Gaussian perturbation. Under regularity conditions on  $\mathbf{X}$  (e.g., finite differential entropy) the perturbation can be made negligible by taking the large  $s$  limit. Finally, let  $\mathbf{X}^* \sim \mathcal{N}(0, \frac{1}{n}I_n)$  be an i.i.d. Gaussian vector with same power as  $\mathbf{X}$ .

An important property of the i.i.d. Gaussian distribution is that, for any  $A \in \mathcal{V}_k(\mathbb{R}^n)$ , the  $k$ -dimensional projection  $A\mathbf{X}^*$  is equal in distribution to the first  $k$  entries in  $\mathbf{X}^*$ . Consequently, the mutual information function satisfies

$$I_{\mathbf{X}^*}(sA^T A) = \frac{k}{n} I_{\mathbf{X}^*}(sI). \quad (26)$$

In the special case where the entries of  $\mathbf{X}$  are independent, it follows from the linear entropy power inequality of Zamir and Feder [19] that

$$I_{\mathbf{X}}(sA^T A) \geq \sum_{i=1}^n [A^T A]_{i,i} I_{X_i}(s). \quad (27)$$

Furthermore, if  $\mathbf{A}$  is drawn uniformly at random from  $\mathcal{V}_k(\mathbb{R}^n)$ , then expectation of the right-hand side is given by

$$\mathbb{E} \left[ \sum_{i=1}^n [A^T A]_{i,i} I_{X_i}(s) \right] = \sum_{i=1}^n \frac{k}{n} I_{X_i}(s) = \frac{k}{n} I_{\mathbf{X}}(sI). \quad (28)$$

Motivated by (28), we define the average EPI deficit to be

$$\delta_{\text{EPI}} \triangleq \frac{1}{k} \mathbb{E} [I_{\mathbf{X}}(sA^T A)] - \frac{1}{n} I_{\mathbf{X}}(sI), \quad (29)$$

where  $\mathbf{A}$  is drawn uniformly at random from  $\mathcal{V}_k(\mathbb{R}^n)$ . If the entries of  $\mathbf{X}$  are not independent, then (27) may not hold. However, as we will see below,  $\delta_{\text{EPI}}$  is non-negative for every distribution on  $\mathbf{X}$ .

Next, we consider the relative entropy between the  $k$ -dimensional projections. For each  $A \in \mathcal{V}_k(\mathbb{R}^n)$ , we have

$$D(P_{\sqrt{s}A\mathbf{X}+\mathbf{N}} \| P_{\sqrt{s}A\mathbf{X}^*+\mathbf{N}}) = \frac{k}{n} I_{\mathbf{X}^*}(sI) - I_{\mathbf{X}}(sA^T A).$$

We define the average CCLT deficit to be

$$\delta_{\text{CCLT}} \triangleq \frac{1}{n} I_{\mathbf{X}^*}(sI) - \frac{1}{k} \mathbb{E} [I_{\mathbf{X}}(sA^T A)], \quad (30)$$

where  $\mathbf{A}$  is drawn uniformly at random from  $\mathcal{V}_k(\mathbb{R}^n)$ .

The following result shows that the average EPI deficit and the average CCLT deficit can both be viewed as measures of the distance between the distribution of  $\mathbf{X}$  and the i.i.d. approximation  $\mathbf{X}^*$ , and that their sum is proportional to the relative entropy after convolution with a Gaussian.

**Theorem 3.** *Let  $\mathbf{X}$  be an  $n$ -dimensional random vector with finite covariance. For all  $s \in (0, \infty)$  and  $k \in \{1, \dots, n-1\}$ ,*

$$\frac{1}{n} D(P_{\sqrt{s}\mathbf{X}+\mathbf{N}} \| P_{\sqrt{s}\mathbf{X}^*+\mathbf{N}}) = \delta_{\text{CCLT}} + \delta_{\text{EPI}}. \quad (31)$$

Furthermore, each of the terms on the right-hand side is non-negative and equal to zero if and only if  $\mathbf{X}$  is i.i.d. Gaussian.

**Remark 4.** Part of the significance of Theorem 3 is that it provides a direct link between the deficit in the EPI and the relative entropy with respect to the i.i.d. Gaussian distribution. Combining results from [18] with the concavity of the mutual information in Theorem 1, it can be shown that, under mild conditions on  $\mathbf{X}$ , the average CCLT deficit  $\delta_{\text{CCLT}}$  is small whenever  $k \ll n$ .

### C. Additivity of information via free probability

Suppose that one obtains outputs from two independent linear Gaussian channels:

$$\mathbf{Y} = A\mathbf{X} + \mathbf{N}, \quad \mathbf{Z} = B\mathbf{X} + \mathbf{N}', \quad (32)$$

where  $A$  and  $B$  both have  $n$  columns and  $\mathbf{N}$  and  $\mathbf{N}'$  are independent standard Gaussian vectors. A natural question of interest is whether the parallel combination of these channels can be characterized in terms of the individual channels. As a direct consequence of the additivity formula in Lemma 2, we see that this question is directly related to additivity of the corresponding matrix SNRs:

$$I(\mathbf{X}; \mathbf{Y}, \mathbf{Z}) = I_{\mathbf{X}}(A^T A + B^T B). \quad (33)$$

In the special case where the signal  $\mathbf{X}$  is i.i.d. Gaussian, the mutual information depends only on the singular values of the channel matrix, or equivalently the eigenvalues of matrix SNR. Consequently, the combination of the channels is completely characterized by the eigenvalues of the sum  $A^T A + B^T B$ . In general, the eigenvalues of the sum of two matrices cannot be determined based only on the eigenvalues of the individual matrices. However, one of the central results from free probability theory is that eigenvalues of the sum of freely independent matrices can be characterized in the large system limit via free additive convolution; see [20], [21].

**Assumption 1.** For each  $n$ ,  $\mathbf{A}_n$  and  $\mathbf{B}_n$  are independent right-orthogonally invariant random matrices with  $n$  columns. Furthermore, as  $n$  increases to infinity, the (random) empirical spectral distributions of  $\mathbf{A}_n^T \mathbf{A}_n$  and  $\mathbf{B}_n^T \mathbf{B}_n$  converge almost surely to compactly supported probability measures  $\mu$  and  $\nu$ .

**Lemma 6** ([20]). *Under Assumption 1, the empirical spectral distribution of the sum  $\mathbf{A}_n^T \mathbf{A}_n + \mathbf{B}_n^T \mathbf{B}_n$  converges almost surely to the probability measured given by the free additive convolution of  $\mu$  and  $\nu$ , which is denoted by  $\mu \boxplus \nu$ .*

For the general case of non-Gaussian signal priors, the mutual information depends on the right-singular vectors of the channel matrix and thus cannot be determined based only on the singular values of the channel matrix. Interestingly though, the additivity principle seen in the case of Gaussian priors is still applicable when the channel matrices are freely independent.

Given a probability measure  $\mu$  on  $[0, \infty)$  we define

$$\mathcal{I}_n(\mu) = \frac{1}{n} \mathbb{E}[I_{\mathbf{X}}(\mathbf{U}^T \Lambda \mathbf{U})], \quad (34)$$

where  $\mathbf{U}$  distributed uniformly on the group of  $n \times n$  orthogonal matrices and  $\Lambda$  is a diagonal matrix whose entries are i.i.d. according to  $\mu$ .

**Assumption 2.** For each  $n$ ,  $\mathbf{X}_n$  is an  $n$ -dimensional random vector with bounded second moment:  $\frac{1}{n} \mathbb{E}[\|\mathbf{X}\|^2] < B$

The next result shows that normalized mutual information associated with a random orthogonally invariant channel matrix converges to its expectation, and furthermore, that the combination of channels with freely independent matrices is characterized by the free additive convolution.

**Theorem 4.** *Under Assumptions 1 and 2, the following convergence holds almost surely in the limit as  $n \rightarrow \infty$ :*

$$\left| \frac{1}{n} I_{\mathbf{X}_n}(\mathbf{A}_n^T \mathbf{A}_n) - \mathcal{I}_n(\mu) \right| \rightarrow 0 \quad (35)$$

$$\left| \frac{1}{n} I_{\mathbf{X}_n}(\mathbf{B}_n^T \mathbf{B}_n) - \mathcal{I}_n(\nu) \right| \rightarrow 0 \quad (36)$$

$$\left| \frac{1}{n} I_{\mathbf{X}_n}(\mathbf{A}_n^T \mathbf{A}_n + \mathbf{B}_n^T \mathbf{B}_n) - \mathcal{I}_n(\mu \boxplus \nu) \right| \rightarrow 0. \quad (37)$$

*Sketch of Proof.* Using the multivariate I-MMSE relationship (10), we show that  $\mathcal{I}_n(\mu)$  is Lipschitz with a constant that depends only on the second-moment of the signal. Concentration with respect to the eigenvectors is established using further Lipschitz properties and standard concentration of measure arguments for the Haar measure on the Stiefel manifold.  $\square$

The question of whether  $\mathcal{I}_n(\mu)$  converges to a well-defined limit is still open in general. Recent work has proved the existence of the limit in the special case of i.i.d. Gaussian matrices [5], [6]. More generally, analysis based on the replica method from statistical physics has provided postulated single-letter formulas for the limit (34) associated with an arbitrary spectral distribution [22], [23]. Recent work by the first author [8] provides an alternative to the replica method that can be applied to the composition of multiple channels.

An important open problem is proving that the conjectured limit of (34) is correct. The significance of Theorem 4 is that it imposes a functional constraint on (34) for large- $n$  that is closely related to the conjectured limit.

## REFERENCES

- [1] T. Richardson, M. A. Shokrollahi, and R. Urbanke, "Design of capacity-approaching irregular low-density parity-check codes," *IEEE Trans. Inform. Theory*, vol. 47, pp. 619–637, Feb. 2001.
- [2] S. Kudekar, S. Kumar, M. Mondelli, H. D. Pfister, E. Şaşıoğlu, and R. Urbanke, "Reed-Muller codes achieve capacity on erasure channels," *IEEE Trans. Inform. Theory*, vol. 63, no. 7, pp. 4298–4316, 2017.
- [3] Y. Kabashima, T. Wadayama, and T. Tanaka, "A typical reconstruction limit for compressed sensing based on  $\ell_p$ -norm minimization," *J. of Stat. Mech.: Theory and Exper.*, vol. 2009, p. L09003, 2009.
- [4] L. Zdeborová and F. Krzakala, "Statistical physics of inference: Thresholds and algorithms," *Adv. in Phys.*, vol. 65, no. 5, pp. 453–552, 2016.
- [5] G. Reeves and H. D. Pfister, "The replica-symmetric prediction for compressed sensing with Gaussian matrices is exact," in *Proc. IEEE Int. Symp. Inform. Theory*, Barcelona, Spain, Jul. 2016, pp. 665 – 669.
- [6] J. Barbier, M. Dia, N. Macris, and F. Krzakala, "The mutual information in random linear estimation," in *Proc. Annual Allerton Conf. on Commun., Control, and Comp.*, Monticello, IL, 2016.
- [7] J. Barbier, F. Krzakala, N. Macris, L. Miolane, and L. Zdeborová, "Phase transitions, optimal errors and optimality of message-passing in generalized linear models," Aug. 2017, [Online]. Available <https://arxiv.org/abs/1708.03395>.
- [8] G. Reeves, "Additivity of information in multilayer networks via additive Gaussian noise transforms," in *Proc. Annual Allerton Conf. on Commun., Control, and Comp.*, Monticello, IL, 2017.
- [9] M. Lamarca, "Linear precoding for mutual information maximization in MIMO systems," in *Proc. of the Int. Conf. on Wireless Comm. Syst.*, Tuscany, Italy, Sep. 2009.
- [10] D. Guo, S. Shamai, and S. Verdú, "Mutual information and minimum mean-square error in Gaussian channels," *IEEE Trans. Inform. Theory*, vol. 51, no. 4, pp. 1261–1282, 2005.
- [11] V. V. Prelov and S. Verdú, "Second-order asymptotics of mutual information," *IEEE Trans. Inform. Theory*, vol. 50, no. 8, pp. 1567–1580, Aug. 2004.
- [12] M. Payaró, M. Gregori, and D. Palomar, "Yet another entropy power inequality with an application," in *Proc. of the Intern. Conf. on Wireless Commun. and Sign. Process.*, Nanjing, China, Nov. 2011.
- [13] O. Rioul, "Information theoretic proofs of entropy power inequalities," *IEEE Trans. Inform. Theory*, vol. 57, no. 1, pp. 33–55, Jan. 2011.
- [14] R. Bustin, M. Payaró, D. P. Palomar, and S. Shamai, "On mmse crossing properties and implications in parallel vector Gaussian channels," *IEEE Trans. Inform. Theory*, vol. 59, no. 2, pp. 818–844, 2013.
- [15] H. L. Van Trees, *Detection, estimation, and modulation theory*. John Wiley & Sons, 2004.
- [16] A. Dembo, "Information inequalities and uncertainty principles," Department of Statistics, Stanford University, Stanford, CA, Tech. Rep. 75, 1990.
- [17] M. Raginsky and I. Sason, *Concentration of Measure Inequalities in Information Theory, Communications, and Coding*, 2nd ed. Foundations and Trends in Communications and Information Theory, 2014.
- [18] G. Reeves, "Conditional central limit theorems for Gaussian projections," in *Proc. IEEE Int. Symp. Inform. Theory*, Aachen, Germany, Jun. 2017, pp. 3055–3059.
- [19] R. Zamir and M. Feder, "A generalization of the entropy power inequality with applications," *IEEE Trans. Inform. Theory*, vol. 39, no. 5, pp. 1723–1728, Sep. 1993.
- [20] D. Voiculescu, "Limit laws for random matrices and free products," *Inventiones mathematicae*, vol. 104, no. 1, pp. 201–220, 1991.
- [21] D. N. C. Tse, "Multiuser receivers, random matrices and free probability," in *Proc. Annual Allerton Conf. on Commun., Control, and Comp.*, 1999.
- [22] K. Takeda, S. Uda, and Y. Kabashima, "Analysis of CDMA systems that are characterized by eigenvalue spectrum," *Europhysics Letters*, vol. 76, no. 6, pp. 1193–1199, Dec. 2006.
- [23] A. Tulino, G. Caire, S. Verdú, and S. Shamai, "Support recovery with sparsely sampled free random matrices," *IEEE Trans. Inform. Theory*, vol. 59, no. 7, pp. 4243–4271, Jul. 2013.