# Noncoherent MIMO Communication: Grassmannian Constellations and Efficient Detection

Ramy H. Gohary, *Member, IEEE*, and Timothy N. Davidson

*Abstract*—This paper considers the design of both a transmitter and a receiver for noncoherent communication over a frequency-flat, richly scattered multiple-input multiple-output (MIMO) channel. The design is guided by the fact that at high signal-to-noise ratios (SNRs), the ergodic capacity of the channel can be achieved by input signals that are isotropically distributed on the (compact) Grassmann manifold. The first part of the paper considers the design of Grassmannian constellations that mimic the isotropic distribution. A subspace perturbation analysis is used to determine an appropriate metric for the distance between Grassmannian constellation points, and using this metric, greedy, direct and rotation-based techniques for designing constellations are proposed. These techniques offer different tradeoffs between the minimum distance of the constellation and the design complexity. In addition, the rotation-based technique results in constellations that have lower storage requirements and admit a natural "quasi-set-partitioning" binary labeling.

In the second part of the paper, a reduced search suboptimum detector is proposed. The development of this detector relies on the subspace perturbation analysis and exploits the geometric properties of the Grassmann manifold and the isotropic distribution of the constellation points and the noise realizations. The performance of this detector is comparable to that of the maximum likelihood detector, but it requires considerably less computational effort. Finally, in order to assess the performance of a given constellation, an exact expression is provided for the pairwise error probability of the ML detector. In comparison to existing pairwise error probability expressions, the proposed expression is numerically stable and does not require the evaluation of residues at poles with high multiplicities.

*Index Terms*—Grassmann manifold, noncoherent multiple input–multiple-output (MIMO) communication, reduced-search maximum likelihood detection, sphere packing.

## I. INTRODUCTION

IN this paper, we consider communication over a block-fading channel in which the channel coefficients are assumed to remain constant for a block of $T$ channel uses and then to change independently to a new channel realization [1]. If neither the transmitter nor the receiver has any *a priori* knowledge of the channel realization (i.e., no channel state information, CSI, is available), then the communication framework is said to be noncoherent [2]–[4], whereas if the receiver has complete *a priori* CSI, the framework is said to be coherent [3]. While the coherent framework facilitates the application of conventional transmission and reception principles, the noncoherent framework has the advantage that it accounts for the communication resources that would have to be expended to obtain CSI at the receiver. Despite the absence of *a priori* CSI at the receiver, noncoherent communication systems with multiple antennas can provide reliable transmission at high data rates. Indeed, it was shown in [3] and [2] that these multiple input–multiple-output (MIMO) noncoherent systems can achieve a significant fraction of the ergodic channel capacity associated with their coherent counterpart[1]s .

For a MIMO communication system with $M$ transmitter antennas and $N$ receiver antennas operating noncoherently over a richly scattered, frequency-flat, block-fading channel with block length $T$, the generic form of the input signals that enable communication at rates approaching the noncoherent ergodic capacity can be expressed as the product of an isotropically distributed $T \times M$ random unitary matrix and a diagonal $M \times M$ matrix $D$ with real nonnegative entries [3]. While this structure of the input signals is capacity achieving irrespective of the values of the received SNR and the channel coherence time $T$, the distribution of the entries of the diagonal matrix $D$ depends on these two factors. For example, it was shown in [5] that at low SNR only one entry of $D$ is nonzero when the transmitter is active. In contrast, achieving capacity for high SNR scenarios requires the input signals to be in the form of isotropically distributed unitary matrices [2], [4], provided that $T$ satisfies

$$T \geq \min\{M, N\} + N. \tag{1}$$

That is, when $T$ satisfies (1), setting $D$ equal to the identity achieves the high SNR ergodic capacity of the noncoherent channel. By comparing the degrees of freedom supported by the unitary component to those supported by the diagonal matrix, $D$, it was concluded in [6] that even at moderate SNRs most of the information will be carried by the unitary component.

By assuming that the communication system operates in the moderate-to-high SNR region, one can gain insight into the

manner in which the coherence time, $T$, affects the achievable data rate. It was shown in [2] and [7] that for given $M$ and $N$, the capacity of the noncoherent channel approaches that of the coherent one as $T$ grows; from which one can conclude that if $T$ is sufficiently long, the amount of time needed for the receiver to acquire a sufficiently accurate channel model becomes insignificant in comparison with the overall signaling interval. On the other hand, when $T$ does not satisfy (1), noncoherent communication can be rather power inefficient. In particular, in the extreme case of $T = 1$, it was shown in [8] that capacity grows only double logarithmically with SNR. The moderate-to-high SNR assumption also provides some insight into how the number of transmit antennas should be chosen for a given block length. In particular, given $T$ and $N$ for a system that satisfies (1), the number of transmit antennas, $M$, required to attain the maximum number of communication degrees of freedom is [2]

$$M = \min\left\{ \left\lfloor \frac{T}{2} \right\rfloor, N \right\}. \qquad (2)$$

In fact, increasing the number of transmit antennas beyond the value in (2) reduces the number of communication degrees of freedom. It is also interesting to note that choosing $N$ to be greater than $M$ does not increase the degrees of freedom [2]. However, choosing $N > M$ does result in an increase in the capacity by an additive term that is independent of the SNR.

In addition to unitary signaling, [4], [9], training-based schemes have been proposed [2], [7], [10], [11] for noncoherent MIMO communication. These schemes comprise a training phase and a coherent communication phase. During the training phase the transmitter sends pilot symbols which are used by the receiver to estimate the channel. Assuming that the channel estimate obtained during the training phase is sufficiently accurate, the receiver then switches to a coherent mode of operation in which the remaining channel coherence time is used to detect the transmitted information coherently. Although training-based schemes were shown [2] to achieve the maximum number of degrees of freedom available for communication at high SNR, they are still short of attaining the full channel capacity, which involves an SNR independent term. This term can be particularly significant when a large number of receive antennas is employed.

In this paper, we will consider communication at moderate-to-high SNRs over a noncoherent MIMO block fading channel in which the block length $T$ satisfies (1) and the number of transmit antennas $M$ satisfies (2). Our approach will be based on the observation that the fading channel matrix does not change the subspace in which the transmitted signal resides. It merely rotates and scales the bases of this subspace. However, the combined effect of noise and fading results in the perturbation of the signal subspace in a specific manner. Based on this geometric insight, it was shown in [2] that, at high SNR, the information carrying object is a linear subspace. That is, information about the transmitted data is contained in the subspace of the received signal and the particular orientation of the received signal vector within the subspace is "informationless". These observations suggest that, for the noncoherent channel, spectrally efficient signaling at high SNR requires the design

of the *bases* of a set of linear signal subspaces rather than the design of the actual signal values. Although the rotation of bases within the subspace is inconsequential, the scaling associated with those bases provides implicit information about the channel. In Section V we will use this fact to develop a reduced complexity detector.

Guided by the results in [2] and [4], in the first part of this paper we design signal constellations that directly mimic the high SNR capacity achieving isotropic distribution. A fundamental issue [12], [13] that arises in constructing a signal constellation is the metric used to measure distances between different constellation points. In [14] multilevel unitary signal constellations were designed using the Kullback-Leibler distance metric. This appears to be a sensible distance metric when constellation points belong to hyperspheres of different radii; a signaling scheme that suits low-to-moderate SNR operation. Our approach to determine the appropriate distance metric for the constellation differs from related work in that it is based on subspace perturbation analysis. This perturbation analysis suggests that from a rate perspective, an appropriate distance measure is given by the chordal Frobenius norm rather than the commonly used projection Frobenius norm (also known as the chordal Frobenius distance) [12], [13].

Apart from the choice of the underlying metric, there are different approaches for generating 'good.' Grassmannian constellations [15]. These approaches can be classified into 1) algebraic (e.g., [16]) and quasi-algebraic (e.g., [17]) approaches, in which the constellations are synthesized from algebraic constructions; 2) approaches that are based on mapping coherent space–time block codes onto the Grassmann manifold (e.g., [18], [19]); and 3) approaches that use various numerical tools for optimization on the Grassmann manifold [9], [20]. The main advantage of the last approach is that it allows the system designer to exploit all the design degrees of freedom without restricting the constellation to have a specific structure. However, irrespective of the underlying distance metric, the direct optimization approach can be quite cumbersome, because the objective is not differentiable, the Grassmann manifold constraint must be enforced, and for large constellations there are many points to design.

In the first part of this paper, we provide several techniques for designing Grassmannian constellations. The first is a greedy technique in which the constellation is constructed sequentially. This technique differs from the surrogate based approach in [9] and the approach in [20] in both the metric used to quantify the distance between constellation points, and the formulation of the underlying optimization problem. In particular, we exploit the smooth geometry of the Grassmann manifold to synthesize an analytic cost function that jointly penalizes the chordal Frobenius norm between the new constellation point and the constellation points already designed. We then minimize this cost function by using a derivative-based optimization algorithm that automatically restricts the iterates to the surface of the Grassmann manifold [21]. Although greedy techniques often result in constellations that provide reasonable performance in practice [20], these techniques are relatively coarse and do not guarantee that the resulting constellation are optimum, even in an asymptotic sense. In order to improve on the greedy approach, we propose two methods for jointly designing Grassmannian constellations.

The first method is direct in the sense that it generates the entire constellation at once. Although the constellations generated by this method possess many desirable features, this method becomes computationally unwieldy as the constellation size increases, and it is currently only appropriate for constellations of up to the order of $2^{10}$ points. In order to reduce the computational complexity, we propose a two-phase design strategy. In the first phase a relatively small Grassmannian constellation (which we will refer to as a proto-constellation) is designed using the direct technique, and in the second phase square unitary matrices are designed to rotate the proto-constellation in an appropriate sense. We will show that rotation preserves the distance between the points of the rotated proto-constellation, and that the performance of constellations designed using the rotation technique lies between that of the constellations designed using the direct and the greedy techniques. Similar to our greedy technique, both the direct and the rotation-based techniques exploit the smooth geometry of the Grassmann manifold and use derivative-based techniques to minimize a smooth cost function that approximates the original nondifferentiable objective. An important advantage of the rotation-based technique over greedy and direct techniques lies in the efficiency with which the constellations can be designed and stored. Furthermore, the inherent structure possessed by rotation-based constellations automatically admits a binary labeling technique that adheres, to a large extent, to the principles that underlie the conventional set-partitioning technique.

In the second part of the paper, we consider the detection and performance analysis of Grassmannian constellations. In particular, we use the subspace perturbation analysis from the first part to develop a reduced-search suboptimum detector. This detector is based on exploiting an inherent property of isotropically distributed Grassmannian constellations and on identifying the role of each component of the received signal matrix. The essence of this detector is to avoid the exhaustive search required for maximum likelihood (ML) detection. Specifically, this detector generates quasi ML decisions by examining a certain subset of the constellation points against the likelihood metric. We demonstrate that, in comparison to ML detection, this detection strategy offers valuable computational savings without significant loss in performance.

Having considered both the design and the efficient detection of Grassmannian constellations, we then consider the analysis of the performance of these constellations at high SNRs. Seeing as the computation of the error probability appears to be mathematically intractable in general, we derive an exact expression for the pairwise error probability (PEP) (e.g., [4], [20]), which can be used to provide further insight into the structure of different Grassmannian constellations. We show that the pairwise error probability corresponds to a specific point on the cumulative distribution function (cdf) of an indefinite quadratic form. As a by-product of our method, we obtain an analytic expression for the whole cdf. This expression can be useful in implementing detection strategies that are based on soft decisions. Our expression for the PEP does not involve the computation of residues required by those in [4] and [20], and hence can be used to study the key parameters that govern the performance of a given con-

stellation; e.g., analyzing the impact of the distance spectrum of the constellation on performance.

The paper is organized as follows. In Section II, we introduce the signal model under consideration. Based on this model, Section III discusses the choice of an appropriate distance metric for the constellation. In Section IV, we propose our constellation design techniques. In Section V, we introduce our reduced search detection strategy and in Section VI we derive an exact expression for the pairwise error probability. Section VII includes our numerical results and Section VIII concludes the paper. For convenience, most of the proofs of our results are deferred to the appendices.

We will adopt the standard notations: $(\cdot)^\dagger$ to denote the Hermitian transpose, $\mathrm{Tr}(\cdot)$ to denote the trace operator of a matrix, $I_q$ to denote the $q \times q$ identity matrix and $I(X; Y)$ to denote the mutual information between $X$ and $Y$. The space of $p \times q$ unitary matrices, where $p \geq q$, will be denoted by $\mathbb{V}_{p,q}$. That is, $\mathbb{V}_{p,q} = \{ Z \in \mathbb{C}^{p \times q} | Z^\dagger Z = I_q \}$. We will also use $[Q]$ to denote the equivalence class represented by $Q \in \mathbb{V}_{T,M}$. That is, $[Q] = \{ \Phi \in \mathbb{V}_{T,M} | \Phi = QP, P \in \mathbb{V}_{M,M} \}$. When $[\cdot]$ is implicit from the context, it will be dropped for simplicity.

## II. System Model

We consider the scenario of moderate-to-high SNR noncoherent communication over a richly scattered frequency-flat block-fading channel with block length $T$ that satisfies (1), $M$ transmit antennas and $N$ receive antennas, where $M$ and $N$ satisfy (2). The transmitter excites the channel in blocks of $T$ channel uses with the rows of the $T \times M$ matrix $Q_X$. The $T \times N$ received signal matrix $Y$ is given by

$$Y = Q_X H + \sqrt{\frac{M}{\rho T}} V \qquad (3)$$

where $H$ is an $M \times N$ channel matrix whose entries are drawn independently from the standard complex Gaussian distribution $\mathcal{CN}(0,1)$, and the $T \times N$ matrix $V$ represents the additive noise whose entries are also drawn from $\mathcal{CN}(0,1)$. The signal-to-noise ratio (SNR) is given by $\rho$ and is independent of the number of transmit antennas $M$.

The capacity achieving input signals at high SNR can be represented by isotropically distributed $M$-dimensional linear subspaces that reside in a larger ambient $T$-dimensional complex Euclidean space, $\mathbb{C}^T$. Since an $M$-dimensional linear subspace of $\mathbb{C}^T$ can be represented by a "tall" $T \times M$ unitary matrix whose columns form a basis for this subspace, we will henceforth assume the input signal matrix $Q_X$ to be unitary; i.e., $Q_X^\dagger Q_X = I_M$. Each of these $M$-dimensional linear subspaces can be regarded as a single point on the compact Grassmann manifold $\mathbb{G}_M(\mathbb{C}^T)$. Since a linear subspace can be specified by an arbitrary basis, points on $\mathbb{G}_M(\mathbb{C}^T)$ are equivalence classes of $T \times M$ unitary matrices, where two matrices are equivalent if they span the same $M$-dimensional subspace. Therefore, the Grassmann manifold can be expressed as [21], $\mathbb{G}_M(\mathbb{C}^T) = \{ [Q] | Q \in \mathbb{V}_{T,M} \}$.

When the signal matrix $Q_X$ in (3) is right multiplied by the $M \times N$ channel matrix $H$, the basis vectors that span the $M$-di-

mensional subspace are rotated and scaled within the same subspace. With this observation [2], one concludes that when the receiver does not know the channel, the particular rotation of the subspace basis is not detectable while the $M$-dimensional linear subspace spanned by this basis is detectable. It follows that the transmitter design problem is to assign information bits to distant linear subspaces, where the distance should be measured in an appropriate sense. The corresponding role of the receiver is to decide on which subspace was transmitted irrespective of its basis. The received signal $Y$, spans an $N$-dimensional subspace. This subspace can be exposed by performing the QR decomposition (in the sense of [22]) on the received signal $Y$; see Section III. In that case, one obtains $Y = Q_Y R_Y$, where $Q_Y$ contains a basis of the $N$-dimensional subspace with some arbitrary orientation and $R_Y$ specifies the scaling and rotation within the subspace. It will be shown in Section V that $R_Y$ is statistically independent of $Q_X$, and that $Q_Y$ is statistically independent of $H$. This indicates that the available information about the channel $H$ appears in $R_Y$, whereas all the information about $Q_X$ is captured by the unitary component $Q_Y$. That is, $Q_Y$ represents the perturbed version of the transmitted subspace available at the receiver. Seeing as the receiver does not have a model for the channel, it attempts to detect the subspace spanned by the columns of $Q_X$ from the subspace spanned by the columns of $Q_Y$ using the implicit channel information contained in $R_Y$.

## III. METRIC CHOICE

In order to design a signal constellation, one needs to define a suitable metric to measure the distance between constellation points. The choice of an appropriate metric is crucial in determining the number of spheres of a given radius that can be packed in a manifold of a specific volume [12], [13]. In the current setting, the manifold under investigation is a Grassmann manifold whose volume is a function of the system SNR. If the specified distance metric does not conform to the underlying communication model, signal points can seem closer or further from each other than they actually are. This may result in system performance measures that do not scale with the increase in SNR in the right sense. (As an example of such scaling, the high SNR ergodic capacity expression in [2] suggests that increasing the data rate by $M(1 - M/T)$ bits per channel use requires an increase of 3 dB in the SNR.) In this section, we study the way in which the noise and the channel perturb the signal subspace. Based on this analysis, we will conclude that the chordal Frobenius norm is an appropriate metric for the rate-centric design of Grassmannian constellations. We will then discuss a number of other distance metrics that have previously been employed in constellation design for the noncoherent MIMO channel.

### A. An Appropriate Distance Metric

Our approach to answering the question of the appropriate distance metric between signal subspaces is based on analysis of the received signal $Y$ in (3). We have stated in Section II that in the absence of noise, the signal subspace is not affected by propagation through the channel [2]. However, when the received signal is contaminated by additive noise, the signal subspace is perturbed in a particular fashion that depends on both the channel and noise. Our goal in this subsection is to identify the effect of different components of the noise and their interaction with the signal subspace.

We begin our analysis by stating a result due to Stewart [22] which we have generalized to the case of $M \leq N$. This result allows us to understand the statistical dependencies between the transmitted and received signal subspaces. We will exploit these dependencies in Section V to develop an efficient detection scheme.

*Lemma 1:* Let $A = QR$, where $A \in \mathbb{C}^{T \times N}$ and $R \in \mathbb{C}^{M \times N}$ have rank $M$ with $M \leq N$, and let $Q \in \mathbb{V}_{T,M}$. Denote the orthogonal complement of $Q$ by $Q^\perp = [Q_1^\perp \quad Q_2^\perp]$, where $Q_1^\perp \in \mathbb{V}_{T,N-M}$ and $Q_2^\perp \in \mathbb{V}_{T,T-N}$. If $E \in \mathbb{C}^{T \times N}$, is such that

$$\tilde{R} + \tilde{Q}^\dagger E \text{ is non- singular} \tag{4}$$

where $\tilde{Q} = [Q \quad Q_1^\perp]$ and $\tilde{R} = [R^\dagger \quad 0_{N \times (N-M)}]^\dagger$, then there exist matrices $W \in \mathbb{C}^{T \times N}$ and $F \in \mathbb{C}^{N \times N}$ that are unique up to right multiplication by $N \times N$ unitary matrix, such that

$$A + E = (\tilde{Q} + W)(\tilde{R} + F) \text{ where } (\tilde{Q} + W) \in \mathbb{V}_{T,N} \tag{5}$$

$$L = (Q_2^\perp)^\dagger E(\tilde{R} + \tilde{Q}^\dagger E)^{-1} \tag{6}$$

$$F = (I + L^\dagger L)^{1/2}(\tilde{R} + \tilde{Q}^\dagger E) - \tilde{R}, \tag{7}$$

$$W = (E - \tilde{Q}F)(\tilde{R} + F)^{-1}. \tag{8}$$

Observe that no structure is imposed on the $R$-matrix in the $QR$ decomposition in Lemma 1, and hence the decomposition is not unique. That is, for any unitary $M \times M$ matrix $P$, $QR$ and $Q'R' = (QP)(P^\dagger R)$ are equivalent $QR$ decompositions of $A$ in the sense of [22]. Equation (5) says that the additive perturbation $E$ results in the augmentation of the $M$-dimensional subspace spanned by the columns of $Q$, namely $[Q] \in \mathbb{G}_M(\mathbb{C}^T)$, to the $N$-dimensional subspace spanned by the columns of $\tilde{Q}+W$, namely $[\tilde{Q}+W] \in \mathbb{G}_N(\mathbb{C}^T)$. Notice that when $N = M$, $\tilde{Q}$ will be equal to $Q$ and the subspace represented by the point $Q$ on $\mathbb{G}_M(\mathbb{C}^T)$ is perturbed to the subspace represented by $Q + W$ on the same manifold.

To apply Lemma 1 to the signal model in (3) we observe that the signal subspace represented by $Q_X$ corresponds to $Q$, and the channel matrix $H$ corresponds to $R$. That is, $Q_X H$ in (3) corresponds to the QR decomposition of $A$ in Lemma 1 and the additive Gaussian noise term $\sqrt{\frac{M}{\rho T}} V$ corresponds to the perturbation matrix $E$. Our goal is to analyze the properties of the perturbation $W$, but before we do so, we need to study the roles played by the different components of the noise term. This will enable us to bound the probability that the condition (4) for Lemma 1 to hold is violated. For the time being we assume that (4) is satisfied. Hence, based on Lemma 1 we can identify the following noise components:

$$G = Q_X^\dagger V \text{ and } \hat{G} = [\hat{G}_1^\dagger \quad \hat{G}_2^\dagger]^\dagger = (Q_X^\perp)^\dagger V \tag{9}$$

where $\hat{G}_1 \in \mathbb{C}^{(N-M) \times N}$ and $\hat{G}_2 \in \mathbb{C}^{(T-N) \times N}$. As will become clear in Section V, each of those components affects the transmitted signal in a certain way. In particular, by observing that $Q_X Q_X^\dagger + Q_{X_1}^\perp (Q_{X_1}^\perp)^\dagger + Q_{X_2}^\perp (Q_{X_2}^\perp)^\dagger = I_T$, where the

partition of $Q_X^\perp$ is conformal with that in Lemma 1, one can rewrite the signal model in (3) as

$$
\begin{aligned}
Y &= Q_X H + \sqrt{\frac{M}{\rho T}} \left( Q_X G + Q_{X_1}^\perp \hat{G}_1 + Q_{X_2}^\perp \hat{G}_2 \right) \\
&= \begin{bmatrix} Q_X & Q_{X_1}^\perp \end{bmatrix} \begin{bmatrix} H + \sqrt{\frac{M}{\rho T}} G \\ \sqrt{\frac{M}{\rho T}} \hat{G}_1 \end{bmatrix} + \sqrt{\frac{M}{\rho T}} Q_{X_2}^\perp \hat{G}_2. \quad (10)
\end{aligned}
$$

Using the decomposition of noise in (10), we observe that

- $G = Q_X^\dagger V$ is a noise component that does not affect the signal subspace. In fact, this noise component contributes to the received signal power; cf. (11) below.
- For $M < N$, the channel matrix is "fat". Hence the product of the $T \times M$ signal $Q_X$ and the $M \times N$ channel matrix $H$ results in the immersion of the $M$-dimensional signal subspace in an $N$-dimensional subspace. The noise component $\hat{G}_1$ spans the range space of $Q_{X_1}^\perp$. If we assume, for the moment, that the noise component $\hat{G}_2$ is zero, then the corresponding received signal $Y$ will be given by the first term on the right hand side of (10). The subspace spanned by the columns of this signal is $N$-dimensional of which only $M$-dimensions are spanned by the signal (and the noise component $G$) and the remaining $(N - M)$ dimensions are spanned by the noise component $\hat{G}_1$. Hence, even in the absence of $\hat{G}_2$, for the receiver to detect the transmitted signal, it will have to decide on which subspace was transmitted from all possible $M$-dimensional subspaces which are spanned by the columns of $Y$. (There are $\binom{N}{M}$ of these $M$-dimensional subspaces.) That is, the noise component $\hat{G}_1$ will introduce ambiguity in distinguishing the signal subspace from the noise subspace. However, it does not change the 'orientation' of the original signal subspace.

The above observations suggest that the perturbations in the signal subspace are due only to the noise component $\hat{G}_2$. We will use this fact later in this section.

We now investigate the applicability of Lemma 1 in the stochastic framework of (3). To that end, we consider the probability that (4) is violated. Let

$$
B = \left[ \left( H + \sqrt{\frac{M}{\rho T}} G \right)^\dagger \quad \sqrt{\frac{M}{\rho T}} \hat{G}_1^\dagger \right]^\dagger \quad (11)
$$

which corresponds to $\tilde{R} + \hat{Q}^\dagger E$ in Lemma 1. In order to bound the probability that (4) is violated, we need to bound the probability of the event that the minimum eigenvalue of $BB^\dagger$, $\lambda_{\min}(BB^\dagger)$, falls below some threshold $\epsilon > 0$. The following lemma, which is based on a result from [23], provides the required bound.

*Lemma 2:* Let $P_v(\epsilon)$ denote the probability that $\lambda_{\min}(BB^\dagger)$ is smaller than some $\epsilon > 0$.

- For $N = M$, $P_v(\epsilon) = 1 - e^{-\rho T M \epsilon / (2\rho T + 2M)} = O\left(\frac{\rho T M \epsilon}{2(\rho T + M)}\right)$.
- For $N > M$, $P_v(\epsilon) \leq 1 - e^{-\rho T N \epsilon / 2M} = O\left(\frac{\rho T N \epsilon}{2M}\right)$.
  *Proof:* See Appendix A.                           □

This lemma confirms that one can use Lemma 1 to investigate signal subspace perturbation for any finite SNR, $\rho$. That is, for $\rho < \infty$, the probability that Lemma 1 does not hold approaches

zero as $\epsilon \to 0$. Moreover, as $\rho \to \infty$, one can see that the perturbation due to the noise term will vanish and the signal will only undergo fading which does not affect the signal subspace.

Lemma 1 states that the subspace $\tilde{Q}$ is perturbed additively by $W$; cf. (8). We now study the properties of $W$ in more detail. We begin by showing that in the general case $W$ belongs to neither the compact Grassmann manifold nor to its tangent space at $Q$. The following lemma exposes the inherent structure of $W$.

*Lemma 3:* Let $U_L \Sigma_L V_L^\dagger$ be the singular value decomposition (SVD) of $L$ defined in (6), and let $\tilde{Q}$ and $Q_2^\perp$ be defined as in Lemma 1. Define $\Xi = \arccos\left( \left( I_N + \Sigma_L^2 \right)^{-1/2} \right)$. Then $W$ in (8) can be written as

$$
W = \begin{bmatrix} \tilde{Q} & Q_2^\perp \end{bmatrix} \begin{bmatrix} -2V_L \sin\left(\frac{1}{2}\Xi\right) \\ 2U_L \cos\left(\frac{1}{2}\Xi\right) \end{bmatrix} \sin\left(\frac{1}{2}\Xi\right) V_L^\dagger. \quad (12)
$$

*Proof:* See Appendix B.                           □

Using the form of the perturbation $W$ in (12), one can verify that,

$$
W^\dagger W = 4V_L \sin^2\left(\frac{1}{2}\Xi\right) V_L^\dagger. \quad (13)
$$

Notice that because $\tilde{R} + F$ in Lemma 1 does not have a particular structure, $W$ in Lemma 1 (and hence in Lemma 3) is unique up to right multiplication by some $N \times N$ unitary matrix. Thus for $W$ to belong to the compact Grassmann manifold, $W^\dagger W$ must be equal to the identity. Simple computation shows that this is the same as requiring $L$ to be equal to a particular multiple of a unitary matrix. Due to the random nature of

$$
L = \hat{G}_2 B^{-1} \quad (14)
$$

this condition is generally not satisfied, indicating that generally $W$ is not a point on the Grassmann manifold.

In order to check whether $W \in \mathbf{T}_Q \mathbb{G}_N(\mathbb{C}^T)$, where $\mathbf{T}_Q \mathbb{G}_{T,N}(\mathbb{C})$ is the tangent space of the Grassmann manifold at $Q$, we use the following fact from [21].

*Lemma 4:* At any point $Q \in \mathbb{G}_N(\mathbb{C}^T)$, the tangent vectors $\Delta \in \mathbf{T}_Q \mathbb{G}_N(\mathbb{C}^T)$ satisfy,

$$
Q^\dagger \Delta = 0. \quad (15)
$$

                                                                          □

Using (12), one can verify that

$$
Q^\dagger W = -2V_L \sin^2\left(\frac{1}{2}\Xi\right) V_L^\dagger. \quad (16)
$$

Hence, for $W$ to belong to $\mathbf{T}_Q \mathbb{G}_N(\mathbb{C}^T)$, $\sin(\frac{1}{2}\Xi)$ must be equal to zero, or, equivalently, $\Sigma_L = 0$. Under that condition, the component of noise in the range space of $Q_2^\perp$ is zero. Given the isotropic nature of noise, the probability associated with that event is vanishingly small.

In order to gain further insight into the nature of the perturbation in Lemma 1, we restrict ourselves to the case of $N = M$. In this case, the model can be simplified by observing that the noise component $\hat{G}_1$ does not exist, and that $\tilde{Q} = Q$ and $Q_2^\perp = Q^\perp$. In Fig. 1 we provide a pictorial representation of the subspace perturbation. The ambient space is the Euclidean space $\mathbb{C}^{T \times M}$.
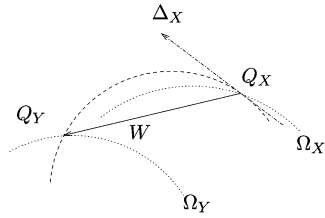
Fig. 1. A pictorial representation of the subspace perturbation analysis: The dashed line represents the geodesic on the Grassmann manifold that passes through the two (dash-dotted) orbits $\Omega_X$ and $\Omega_Y$. The vector $\Delta_X$ lies in $\mathbf{T}_{Q_X}\mathbb{G}_N(\mathbb{C}^T)$, the tangent space at $Q_X$, and the vector $W$ denotes the perturbation.

The points $Q_X$ and $Q_Y$ lie on the Grassmannian manifold and the subspaces that their columns span are represented by the orbits $\Omega_X$ and $\Omega_Y$ in $\mathbb{C}^{T \times M}$.[1] The two orbits are connected by a geodesic on the Grassmannian manifold (depicted by the dashed line), whose length gives the geodesic distance between $Q_X$ and $Q_Y$. The vector $\Delta_X$ is the tangent to the geodesic at $Q_X$. As illustrated in Fig. 1, the perturbation, $W$, of $Q_X$ to $Q_Y$ is additive in the Euclidean space $\mathbb{C}^{T \times M}$ and does not lie along the geodesic nor does it lie in the tangent space. As suggested by Fig. 1, the norm of $W$ quantifies the perturbation of the subspace spanned by $Q_X$ (i.e., $\Omega_X$) to that spanned by $Q_Y$ (i.e., $\Omega_Y$). In the following theorem, we formally state that result. (The tangential and normal components of $W$ are $Q_X^\dagger W$ and $(Q_X^\perp)^\dagger W$, respectively [21], and the chordal Frobenius norm is defined in (19), below.)

*Theorem 1:* Let $N = M$, and let $Y = Q_Y R_Y = Q_X H + V$, where $Q_X, Q_Y \in \mathbb{V}_{T,M}$, $H \in \mathbb{C}^{M \times M}$ and $V \in \mathbb{C}^{T \times M}$. Let the subspace spanned by the columns of $Q_Y$ be $\Omega_Y$ and that spanned by the columns of $Q_X$ be $\Omega_X$. Then the basis of $\Omega_Y$ is an additively perturbed version of the basis of $\Omega_X$ and the Frobenius norm of the perturbation vector is given by the chordal Frobenius norm between $\Omega_X$ and $\Omega_Y$. Moreover, the tangential and normal components of the perturbation $W$ and hence the norm of the perturbation vector are statistically independent of $Q_X$.

*Proof:* See Appendix C. □

*Remark 1:* From Lemma 3, one can see that while the norm of $W$ does not depend on $Q_X$, $W$ itself depends on $Q_X$. This is in contrast to typical coherent communication scenarios in which the perturbation induced by noise does not depend on the transmitted signal. To further investigate this dependence, let us consider a set of noise and channel realizations indexed by $t$, namely $\{(H_t, V_t)\}$, and let $\{W_X(H_t, V_t)\}$ denote the set of perturbations that these channel and noise realizations induce on the point $Q_X$. Since each of the resulting unitary components of the received signal, $Q_{Y,t}$, satisfies $Q_{Y,t}^\dagger Q_{Y,t} = I$, each perturbation satisfies

$$(Q_X + W_X(H_t, V_t))^\dagger (Q_X + W_X(H_t, V_t)) = I. \quad (17)$$

[1] The orbit of a point in the Grassmann manifold is the set of points in the Euclidean space that span the same subspace. This orbit can be generated by the right action of the group of $M \times M$ unitary matrices on any "tall" $T \times M$ unitary matrix that represents the point on the Grassmannian manifold.
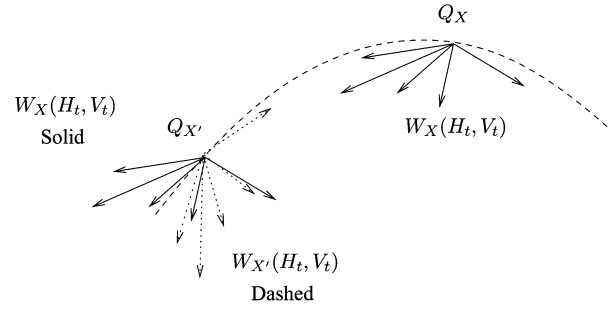


Fig. 2. A pictorial representation of the dependence of the perturbation $W$ on $Q_X$. The dashed curve joining the points $Q_X$ and $Q_{X'}$ on the Grassmannian manifold represents the geodesic between those points. For a given set of channel and noise realizations, $\{(H_t, V_t)\}$, we have illustrated the perturbations induced on $Q_X$ and $Q_{X'}$, namely $\{W_X(H_t, V_t)\}$ (solid rays) and $\{W_{X'}(H_t, V_t)\}$ (dotted rays), respectively. For reference, we have parallel translated $\{W_X(H_t, V_t)\}$ to the point $Q_{X'}$, and the figure illustrates the fact that this is not the set of perturbations that these channel and noise realizations induce on $Q_{X'}$.

If we were to parallel translate $\{W_X(H_t, V_t)\}$ to another point on the manifold, say, $Q_{X'}$, then it is clear that $Q_{X'} + W_X(H_t, V_t)$ will not necessarily satisfy a relation of the form in (17), and hence that $Q_{X'} + W_X(H_t, V_t)$ is not necessarily on the manifold $\mathbb{G}_M(\mathbb{C}^T)$; see Fig. 2 for a pictorial representation. Therefore, it is clear that $W_X(H_t, V_t)$ cannot be independent of $Q_X$. □

Having taken a closer look at the way in which the noise perturbs the signal subspace, we are now in a position to design signal constellations that enable the high SNR ergodic capacity of the noncoherent MIMO channel to be approached. In order to do that, we recall that achieving this capacity corresponds to packing spheres in the Grassmann manifold, where the center of each sphere represents a constellation point and the radius corresponds to the region within which the noise perturbs the constellation point with high probability; cf. [24]. We have seen in Lemma 1 that the combined effect of the channel and the noise is to perturb the signal subspace spanned by the columns of $Q_X$ to another subspace which is spanned by the columns of $Q_Y = \tilde{Q}_X + W \in \mathbb{G}_N(\mathbb{C}^T)$, where $\tilde{Q}_X$ was defined in that lemma, and we have seen in Theorem 1 that for $N = M$ the norm of this perturbation is given by the chordal Frobenius norm between $Q_X$ and $Q_Y$. These observations suggest that the chordal Frobenius norm is an appropriate distance metric for the rate-centric design of Grassmannian constellations, and we will use this metric in our designs.

We will formally define the chordal Frobenius norm below, but first we state a result from [25].

*Lemma 5:* Let $A \in \mathbb{C}^{M \times M}$ and let $A = U_A \Sigma_A V_A^\dagger$ denote its SVD. Then

$$\max_{Z \in \mathbb{C}^{M \times M}, ZZ^\dagger = I} \text{Tr}(Z^\dagger A^\dagger) = \max_{Z \in \mathbb{C}^{M \times M}, ZZ^\dagger = I} \text{Tr}(AZ)$$
$$= \text{Tr}(\Sigma_A) \quad (18)$$

and the optimal $Z$ is given by $Z = V_A U_A^\dagger$. □

The chordal Frobenius norm between the two subspaces that are spanned by the columns of $Q_{Y_1}$ and $Q_{Y_2}$ is defined [21] to be the (nonnegative) square root of $d^2(Q_{Y_1}, Q_{Y_2})$, where

$d^2(Q_{Y_1}, Q_{Y_2})$ is given in (19) at the bottom of the page, where $Q_{Y_1}^\dagger Q_{Y_2} = U_{Q_{Y_1}^\dagger Q_{Y_2}} \Sigma_{Q_{Y_1}^\dagger Q_{Y_2}} V_{Q_{Y_1}^\dagger Q_{Y_2}}^\dagger$ is the singular value decomposition (SVD) of $Q_{Y_1}^\dagger Q_{Y_2}$. Using Lemma 5, we have

$$d^2(Q_{Y_1}, Q_{Y_2}) = 2M - 2\mathrm{Tr}\left(\Sigma_{Q_{Y_1}^\dagger Q_{Y_2}}\right). \qquad (20)$$

### B. Other Distance Metrics

In the previous subsection we have used subspace perturbation analysis to conclude that from the ergodic rate perspective, the chordal Frobenius norm is an appropriate distance metric for packing spheres in the Grassmann manifold. In this subsection we will discuss the relationships between this metric and other distance metrics.

The discussion leading to Theorem 1 suggests that the combined effect of the channel and noise does not perturb the signal subspace along a geodesic on the Grassmann manifold. Instead, it takes the signal along a more general path that resides in the ambient Euclidean space, $\mathbb{C}^{T \times M}$. (Note that $\mathbb{C}^{T \times M}$ has $2TM$ degrees of freedom whereas $\mathbb{G}_T(\mathbb{C}^M)$ has only $2M(T - M)$.) This observation suggests that the arc length (used in [18]) might not be the appropriate distance metric for constellation design and symbol detection. Another metric that has been used in [9], [26] for designing Grassmannian constellations and packings [26] is the projection Frobenius norm defined as,[2]

$$d_p^2(Q_{Y_1}, Q_{Y_2}) = \frac{1}{2}\|Q_{Y_1}Q_{Y_1}^\dagger - Q_{Y_2}Q_{Y_2}^\dagger\|^2$$
$$= M - \mathrm{Tr}\left(\Sigma_{Q_{Y_1}^\dagger Q_{Y_2}}^2\right). \qquad (21)$$

This distance metric results from embedding the Grassmann manifold in the space of $T \times T$ projection matrices of rank $M$ [12], [13], [21]. However, the norm defined through this embedding is strictly less than the norm defined by observing the Grassmann manifold, $\mathbb{G}_M(\mathbb{C}^T)$, as a subspace of the Euclidean space $\mathbb{C}^{T \times M}$. This is a consequence of the fact that the space of $T \times T$ projection matrices of rank $M$ is of higher dimension than the Euclidean space $\mathbb{C}^{T \times M}$. By moving along higher dimensional paths, we may "cut corners" in measuring the distance between any two points [21]. However, the perturbation analysis in Lemma 1 shows that the noise perturbs the subspace by an additive term ($W$ in (8)), and hence the perturbation path cannot lie in a space of higher dimension than the Euclidean space $\mathbb{C}^{T \times M}$. That said, the projection Frobenius norm is a lower bound on the chordal Frobenius norm, and hence constellations that are designed using the projection Frobenius norm may also exhibit favourable performance characteristics.

[2]This norm was called the chordal Frobenius distance in [26].

A comprehensive discussion on other distance metrics and the corresponding embeddings can be found in [21] and [12]. It is worth mentioning that the impact of choosing the appropriate metric on the density of packing becomes less acute as the SNR increases, because several metrics become equivalent in the limit as the angles between subspaces approach zero. This observation can be easily verified by taking the limits on the bounds in [12] as the radius of the metric balls approaches zero.

*Remark 2:* Although a Grassmannian constellation that corresponds to densely packed spheres in the Grassmann manifold is sufficient to approach the noncoherent ergodic capacity, this constellation does not necessarily achieve maximal diversity; cf. [20, SectionVI-B]. In particular, the radius of the spheres in a capacity-approaching packing is chosen such that the perturbed constellation point lies within the sphere with high probability; cf. [24]. That is, for this packing the radius of spheres is determined by the probability distribution of the perturbation around the origin. (For example, in standard additive white Gaussian channels, this radius is equal to the square root of the noise variance.) In contrast, for a constellation with a given cardinality to achieve maximal diversity, the constellation points are designed in such a way that the high SNR rate of decay of error probability is maximal. Hence, the assumption that underlies such a design is that the constellation is used to operate away from the ergodic capacity. Unlike a capacity-approaching constellation, the spacing between constellation points in a constellation with maximal diversity is determined by the tail of the perturbation distribution. ☐

## IV. CONSTELLATION DESIGNS

Having shown that the chordal Frobenius norm is an appropriate metric for measuring distances between constellation points on the Grassmann manifold, we now develop practical procedure s for the design of Grassmannian constellations. We will adopt a design approach similar to the one in [9], in which the number of spheres is given, and the minimum of the radii of the spheres is to be maximized. In particular, given $|\mathcal{C}|$, the constellation design problem consists of finding unitary matrices $Q_{X_i}$ with maximum pairwise chordal Frobenius norm between the subspaces they span. Given (19), that problem can be formulated as

$$\min_{\{Q_{X_r}\}_{r=1}^{|\mathcal{C}|}} \max_{1 \le i,j \le |\mathcal{C}|} \quad \mathrm{Tr}(\Sigma_{ij})$$
$$\text{subject to} \quad Q_{X_k} \in \mathbb{G}_M(\mathbb{C}^T),$$
$$\forall k \in \{1, 2, \ldots, |\mathcal{C}|\} \qquad (22)$$

where $\Sigma_{ij}$ denotes the $M \times M$ diagonal matrix of singular values of $Q_{X_i}^\dagger Q_{X_j}$.

---

$$d^2(Q_{Y_1}, Q_{Y_2}) = \min_{Z \in \mathbb{C}^{M \times M}, \, ZZ^\dagger = I} \|Q_{Y_1} - Q_{Y_2}Z\|^2$$
$$= 2M - \max_{Z \in \mathbb{C}^{M \times M}, \, ZZ^\dagger = I} \left\{ \mathrm{Tr}\left(Z^\dagger Q_{Y_2}^\dagger Q_{Y_1}\right) + \mathrm{Tr}\left(Q_{Y_1}^\dagger Q_{Y_2}Z\right) \right\} \qquad (19)$$

The development of effective algorithms for (approximately) solving (22) requires the resolution of two key issues: the nondifferentiability of the objective, due to the presence of the $\max(\cdot)$ function, and optimizing over multiple Grassmannian points at the same time. In Sections IV-A–C, we will propose three design techniques that provide different approaches to addressing these two issues. These techniques offer the designer a tradeoff between the minimum mutual distance between constellation points and the design complexity.

### A. A Greedy Algorithm

A relatively coarse method for finding an approximate solution to the constellation design problem in (22) is to generate the constellation sequentially using a greedy algorithm. Starting from an arbitrary point on the Grassmann manifold, the essence of the greedy algorithm is to augment the constellation recursively by one constellation point at a time. Since the isotropic distribution implies maximum distance between constellation points, given the set of current constellation points, we would like to choose the next constellation point to be the one that maximizes the minimum distance to all points in the set. That is, assume that $Q_{X_k}, k \in \{1, 2, \ldots, i-1\}$ have already been determined, then, for $i = 2, \ldots, |\mathcal{C}|$, we choose

$$
\begin{aligned}
Q_{X_i} &= \arg \max_{\{Q|Q^\dagger Q=I\}} \min_{1 \le j \le i-1} d(Q, Q_{X_j}) \\
&= \arg \min_{\{Q|Q^\dagger Q=I\}} \max_{1 \le j \le i-1} \operatorname{Tr}\big(\Sigma_{Q^\dagger Q_{X_j}}\big), \quad i = 2, \ldots, |\mathcal{C}|
\end{aligned}
\tag{23}
$$

where $\Sigma_{Q^\dagger Q_{X_j}}$ is the $M \times M$ diagonal matrix of singular values of $Q^\dagger Q_{X_j}$; see (20). For the special case in which the block length $T$ is an even number and the number of transmit antennas $M$ is chosen to maximize the number of degrees of freedom, that is $M = T/2$, we have the following proposition that can be used to reduce the computational complexity by a factor of two.

*Proposition 1:* For $T = 2M$, if $|\mathcal{C}| = 2K$, for some integer $K$, and if $Q_{X_i} \in \mathcal{C}$ is generated by (23) then

$$
Q_{X_i}^\perp \in \mathcal{C}.
$$

*Proof:* See Appendix D. $\quad\square$

*Remark 3:* The statement of Proposition 1 also holds if the projection Frobenius norm (cf. (21)) is used instead of the chordal Frobenius norm (cf. (20)). This can be proved using an argument similar to the one given in Appendix D. $\quad\square$

When the conditions of Proposition 1 are satisfied, the result in Lemma 8 in Appendix D implies that the design problem in (23) can be written as shown in (24) at the bottom of the page with the remaining points being $\{Q_{X_k}^\perp\}$.

In order to solve (23) (or (24)), one needs to perform the inner maximization over the discrete set of $Q_{X_j}$, $j \in \{1, i-1\}$. This discrete minimization, along with the fact that $\max_{1 \le j \le i-1} \operatorname{Tr}(\Sigma_{Q^\dagger Q_{X_j}})$ is not smooth in $Q$, means that (23) is particularly difficult to solve. Instead of performing explicit maximization, we propose to approximate the $\max(\cdot)$ function in (23) (or (24)) by a smooth differentiable function that is amenable to effective gradient-based numerical optimization techniques. In order to provide such an approximation, we note that the "Jacobian logarithm" (e.g., [27]) for two real numbers $a$ and $b$ is

$$
\log(e^a + e^b) = m\Big(1 + m^{-1}\log(1 + e^{-|a-b|})\Big)
\tag{25}
$$

where $m = \max(a, b)$, and the second term in (25) is a correction term. For sufficiently large values of $|a - b|$, the correction term goes to zero and $\max(a, b)$ can be well approximated by $\log(e^a + e^b)$. To refine this approximation for $a, b > 1$, we propose to use the function

$$
\begin{aligned}
F_n(a, b) &= \Big(\log(e^{a^n} + e^{b^n})\Big)^{1/n} \\
&= m\Big(1 + m^{-n}\log(1 + e^{-|a^n - b^n|})\Big)^{1/n}
\end{aligned}
\tag{26}
$$

for $n \ge 1$ to reduce the effect of the correction term. In fact, for any $a, b > 1$ as $n \to \infty$, $F_n(a, b) \to \max(a, b)$. Using (26), the smooth approximation of the optimization problem in (22) can be written as[3]

$$
Q_{X_i} = \arg \min_{\{Q|Q^\dagger Q=I\}} \Bigg( \log\Big( \sum_{j=1}^{i-1} e^{\operatorname{Tr}^n\big(\Sigma_{Q^\dagger Q_{X_j}}\big)} \Big) \Bigg)^{1/n},
$$
$$
i = 2, \ldots, |\mathcal{C}|. \tag{27}
$$

As mentioned above, $d(Q_Y, Q_{Y_0})$, and hence $\operatorname{Tr}(\Sigma_{Q_Y^\dagger Q_{Y_0}})$ (cf. (20)), are functions on the Grassmann manifold. The Grassmann manifold is a smooth topological space [21], and since the objective function in (27) is also smooth, one can use efficient gradient-based optimization algorithms to solve the optimization problem in (27). A particular class of these algorithms is presented in [21]. The optimization algorithms in that class start from an initial point on the manifold, and subsequent iterates are generated by moving along geodesics of the manifold. This approach results in unconstrained optimization problems in which the required orthogonality constraints are automatically satisfied at each iteration. In our numerical work, we will use the Riemannian version of the conjugate gradient algorithm (with resets) given in [21].

Our greedy algorithm has a number of interesting properties. In particular:

[3]Our numerical experiments have shown that the convergence rate of such algorithms can be substantially improved by (partially) solving a sequence of problems indexed by increasing values of $n$.

$$
Q_{X_i} = \arg \min_{\{Q|Q^\dagger Q=I\}} \max_{1 \le j \le i-1} \max\Bigg\{ \operatorname{Tr}\Big(\Sigma_{Q^\dagger Q_{X_j}}\Big), \operatorname{Tr}\Big(\big(I_M - \Sigma_{Q^\dagger Q_{X_j}}^2\big)^{1/2}\Big) \Bigg\}, \quad i = 2, \ldots, |\mathcal{C}|/2
\tag{24}
$$

- The greedy algorithm involves optimization of one constellation point at a time, and hence each instance of (27) iteration of the greedy algorithm involves optimization over $2M(T - M)$ real dimensions. Therefore, the accuracy of the current solution of the conjugate-gradient algorithm (with resets) on the Grassmann manifold doubles every $2M(T - M)$ iterations; see [21, Sec. 3.5.1]. For practical communication systems, this number is relatively small, and hence each step of the greedy algorithm typically exhibits fast convergence.

- Although the $\max$ function in (23) can be well-approximated by increasing the value of $n$ in (26), the greedy algorithm remains inherently suboptimal due to the underlying sequential design.

- One of the attractive features possessed by the constellations generated by the greedy algorithm is that each constellation is a subset of larger constellations. This feature enables the transmitter to easily vary the transmission rate; something that can be particularly useful if the Grassmannian constellation were to be used in an adaptive coded modulation (ACM) framework [28].

Having provided a greedy technique for designing Grassmannian constellations, we now propose alternative methods that enable these constellations to be jointly designed.

### B. Direct Design

In this method, we directly address the problem of simultaneous design of the $|\mathcal{C}|$ Grassmannian points formulated in (22). Our first step is to apply (26) to obtain the following smooth approximation of (22):

$$\min_{\{Q_{X_r}\}_{r=1}^{|\mathcal{C}|}} \quad \left(\log\left(\sum_{i=1}^{|\mathcal{C}|-1} \sum_{j=i+1}^{|\mathcal{C}|} e^{\mathrm{Tr}^n(\Sigma_{ij})}\right)\right)^{1/n} \quad \text{(28a)}$$

$$\text{subject to} \quad Q_{X_k} \in \mathbb{G}_M(\mathbb{C}^T), \quad k = 1, \ldots, |\mathcal{C}|. \quad \text{(28b)}$$

Notice that, unlike the greedy technique, the constellation points that solve (28) approach $\mathcal{C}^o$ as $n \to \infty$, where $\mathcal{C}^o$ is a Grassmannian constellation with maximal minimum distance. This asymptotic optimality could not be guaranteed for the greedy technique because of the sequential fashion in which that technique generates the constellation points. Although it is smooth, the optimization problem in (28) is over multiple points on the manifold, and hence the techniques from [21] that were employed for the greedy algorithm are not immediately applicable. In order to facilitate the adaptation of such techniques to (28), we construct the block diagonal $|\mathcal{C}|T \times |\mathcal{C}|M$ matrix $\bar{Q}$ with $|\mathcal{C}|$ diagonal blocks of size $T \times M$; i.e.

$$\bar{Q} = \mathrm{blkdiag}(Q_{X_1}, \ldots, Q_{X_{|\mathcal{C}|}}) \quad \text{(29)}$$

where $\mathrm{blkdiag}$ is the block diagonal operator. Each matrix $\Sigma_{ij}$ in (22) and (28) can be expressed as

$$\Sigma_{ij} = U_{ij}^\dagger I_M^{(i)} \bar{Q}^\dagger \left(I_T^{(i)}\right)^\dagger I_T^{(j)} \bar{Q} \left(I_M^{(j)}\right)^\dagger V_{ij} \quad \text{(30)}$$

where $U_{ij}\Sigma_{ij}V_{ij}^\dagger$ denotes the singular value decomposition of $Q_{X_i}^\dagger Q_{X_j}$ and the matrix $I_K^{(\ell)}$ denotes the all zero $K \times |\mathcal{C}|K$ "fat" matrix with the $\ell$th $K \times K$ block replaced with $I_K$.

By using (29) and (30) one can reformulate (28) as an optimization problem over the matrix $\bar{Q}$, which represents the $|\mathcal{C}|$ points in $\mathbb{G}_M^T$ by a single point on the Grassmann manifold $\mathbb{G}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$. (The dimension of $\mathbb{G}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$ is $2|\mathcal{C}|^2 M(T - M)$, whereas that of $\mathbb{G}_M(\mathbb{C})^T$ is $2M(T - M)$.) However, because the matrix $\bar{Q}$ is restricted to have a block diagonal structure, $\bar{Q}$ resides in a submanifold of $\mathbb{G}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$ of dimension $2|\mathcal{C}|M(T - M) = |\mathcal{C}| \times \dim(\mathbb{G}_M(\mathbb{C}^T))$. We will refer to this submanifold as $\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$. In order to adapt the techniques in [21] to this submanifold, we need to examine its tangent space. First, the submanifold inherits both the canonical inner product and the projector [21] from the original $\mathbb{G}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$. Now, for all $\bar{Q} \in \tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$

$$\dim\left(\mathbf{T}_{\bar{Q}}\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})\right) = \dim\left(\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})\right)$$

where $\mathbf{T}_{\bar{Q}}\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$ denotes the tangent space to $\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$ at $\bar{Q}$. Since tangent vectors $\Delta \in \mathbf{T}_{\bar{Q}}\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$ satisfy $\Delta^\dagger \bar{Q} = 0$, (cf. (15)) it is clear that tangent vectors also possess a block diagonal structure. That is, the tangent vectors $\Delta \in \mathbf{T}_{\bar{Q}}\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$ can be expressed as the block diagonal component of $(I_{T|\mathcal{C}|} - \bar{Q}\bar{Q}^\dagger)X$, where $X \in \mathbb{C}^{T|\mathcal{C}| \times M|\mathcal{C}|}$. This implies that if the gradient-based algorithms for optimization on the Grassmann manifold (e.g., [21]) are initialized with a point on the submanifold $\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$, then the iterates are guaranteed to remain on this submanifold, and hence those algorithms can now be extended to the case of direct design.

In order to illustrate the potential impact of the direct design technique, in Fig. 3 we compare the distance spectrum of 16-point constellations designed using the greedy and direct techniques, for a system with $M = 2$ and $T = 4$. For the greedy technique the minimum distance is about 1.0917, whereas for the joint design the minimum distance is about 1.1759. (In order to simplify the design, we enforced the antipodal symmetry of Proposition 1 in both designs.) In order to demonstrate the corresponding performance of these constellations when used in noncoherent communication, in Fig. 4 we plot the block error rate of the maximum likelihood detector for the case in which the number of receive antennas $N = 2$. At a block error rate of $10^{-5}$, the direct constellation design yields an SNR gain of about 6 dB over the greedy technique.

### C. Rotation-Based Design

Although the direct design technique generates good constellations, the size of the matrix $\bar{Q}$ in (29) results in an unwieldy optimization problem for large constellations. Furthermore, the resulting constellation points do not have a specific structure that could be used, among other things, to reduce the memory required for their storage.

Since the Grassmann manifold is a Lie group under left multiplication by square unitary matrices, every element $Q_{X_k} \in \mathbb{G}_M(\mathbb{C}^T)$ can be expressed as a rotation $\Phi$ of another element, say $Q_{X_j}$. That is, $Q_{X_k} = \Phi Q_{X_j}$. Using the block diagonal formulation presented in Section IV-B, it can be readily seen that a constellation of $|\mathcal{C}|$ points in $\mathbb{G}_M(\mathbb{C}^T)$ is equivalent to one point in $\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$. Since $\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$ is a Lie group under left

(a) Greedy, $d_{\min} = 1.0917$
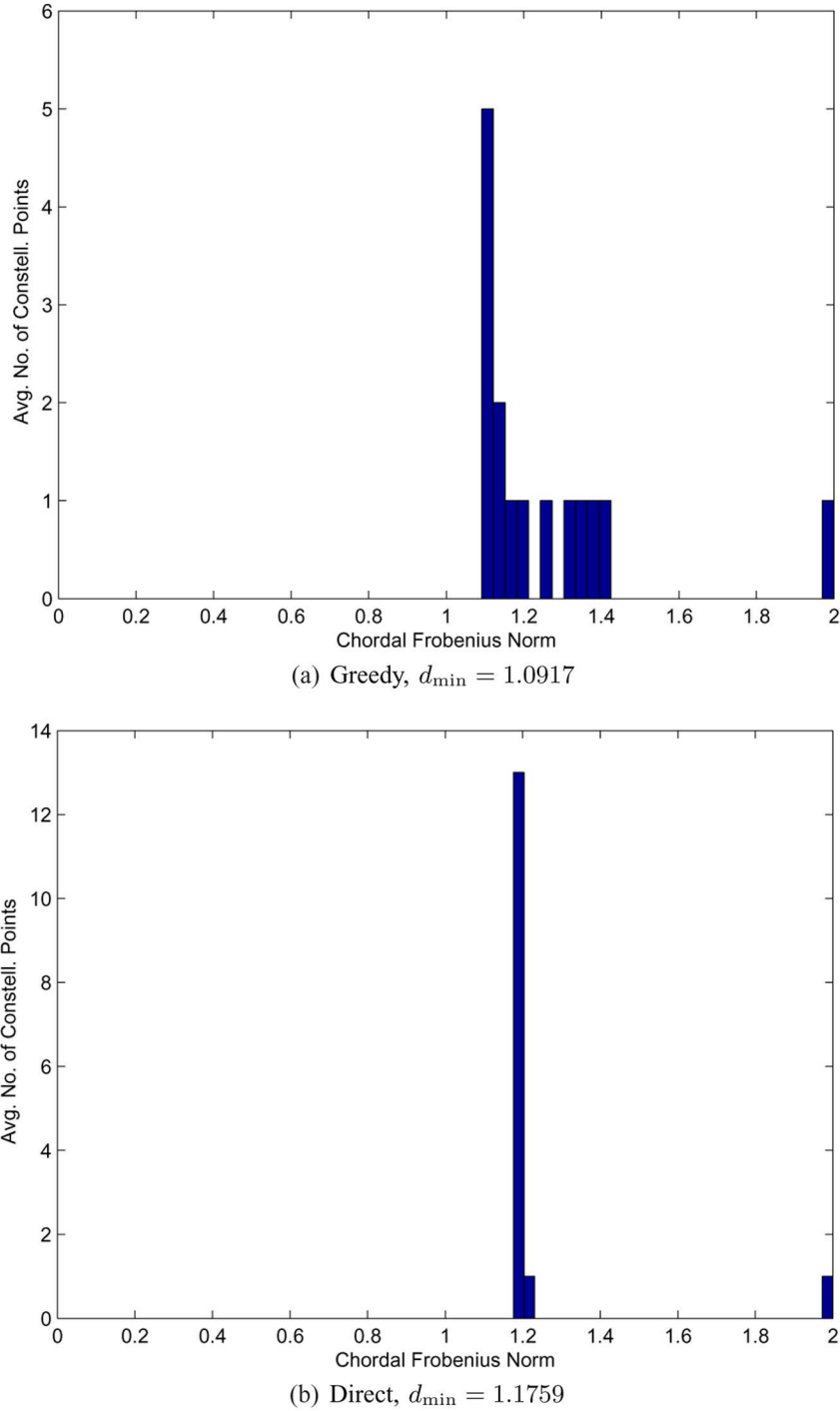


(b) Direct, $d_{\min} = 1.1759$

Fig. 3. Distance spectra of greedy and directly designed 16-point constellations for a system with $M = 2$ and $T = 4$.

multiplication by block diagonal square unitary matrices, a large Grassmannian constellation $\mathcal{C}_L$ of size $L|\mathcal{C}|$ can be represented as a collection of $|\mathcal{C}|$ points in $\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$ and a set of $L$ rotation matrices $\{\bar{\Phi}\}_{i=1}^{L}$. Since $\mathcal{C}$ can be represented as a single point in $\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$, (cf. (29)), $\mathcal{C}_L$ can be represented by $L$ points in $\tilde{\mathbb{G}}_{M|\mathcal{C}|}(\mathbb{C}^{T|\mathcal{C}|})$, where each point is a specific left multiplication of the representation of $\mathcal{C}$. In other words

$$\mathcal{C}_L = \coprod_{i=1}^{L} \bar{\Phi}_i \mathcal{C} \qquad (31)$$

where $\coprod$ denotes the disjoint union operation, and $\mathcal{C}$ and $\mathcal{C}_L$ are represented by one and $L$ block diagonal matrices, respectively.

In order to generate a constellation $\mathcal{C}_L$ with maximal minimum distance, one ought to design both the small constellation $\mathcal{C}$ (which we will call the proto-constellation) and the rotation matrices $\{\bar{\Phi}_i\}$ jointly. However, this problem is even more complicated than designing $\mathcal{C}_L$ directly using (22), because $\bar{\Phi}_i$ represent points on the unitary group $\mathbb{U}_T$, whereas the constellation $\mathcal{C}$ consists of points in $\mathbb{G}_M(\mathbb{C}^T)$. Instead, we propose a simplified two-step approach. The first step is to generate a proto-constel-
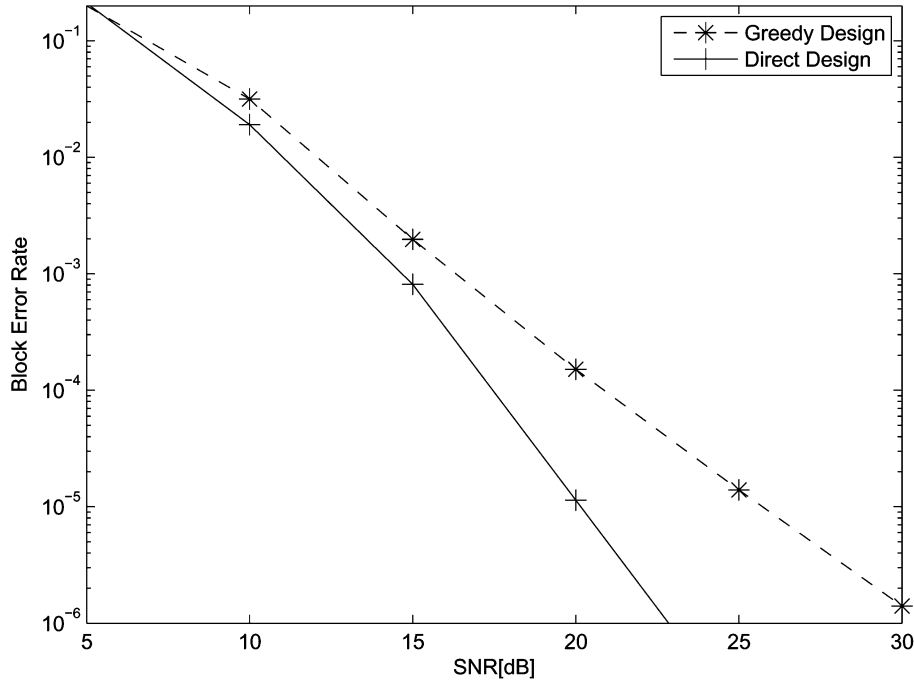
Fig. 4. Performance comparison of the greedy and directly designed constellations in Fig. 3 for a system with $N = 2$ receiver antennas.

lation $\mathcal{C}$ using the direct technique. The second step is to restrict the rotation matrices $\{\bar{\Phi}_i\}$ in (31) to have identical blocks; i.e., $\{\bar{\Phi}_i = I_{|\mathcal{C}|} \otimes \Phi_i\}_{i=1}^{L-1}$. In this case the design of the rotation matrices can be formulated as

$$\min \quad \left( \log\left( \sum_{i,j=1}^{|\mathcal{C}|} \sum_{k=1}^{L-1} \sum_{\ell=k+1}^{L} e^{\text{Tr}^n\left(\Sigma_{ij}^{k\ell}\right)} \right) \right)^{1/n} \quad (32a)$$

$$\text{subject to} \quad \Phi_k \in \mathbb{U}_T, \quad k = 1, \ldots, L-1 \quad (32b)$$

where $\Sigma_{ij}^{k\ell}$ denotes the $M \times M$ the matrices of singular values of $Q_{X_i}^{\dagger} \Phi_k^{\dagger} \Phi_\ell Q_{X_j}$. Like the Grassmannian manifold, the unitary group, $\mathbb{U}_T$, is a smooth manifold to which the gradient-based optimization techniques presented in [21] can be effectively applied. As in the greedy and direct methods, we will use the Riemannian version of the conjugate-gradient algorithm (with resets) in [21].

We now state two properties of the rotated constellations.

*Property 1:* For any two elements $Q_{X_1}, Q_{X_2} \in \mathbb{G}_M(\mathbb{C}^T)$ and any rotation $\Phi$

$$d(Q_{X_1}, Q_{X_2}) = d(\Phi Q_{X_1}, \Phi Q_{X_2}). \quad (33)$$

*Proof:* This property can be verified by evaluating the right hand side of the equality and showing that it is independent of $\Phi$. $\square$

This property guarantees that rotating the proto-constellation by a block diagonal unitary matrix with identical blocks preserves the distance between the points of the rotated proto-constellation. That is, in designing rotations, one only needs to maximize the distance between points that belong to different rotated versions of the proto-constellation.

*Property 2:* For any $Q_X \in \mathbb{G}_M(\mathbb{C}^T)$ and $\Phi \in \mathbb{U}_T$, $\Phi Q_X^{\perp} = (\Phi Q_X)^{\perp}$.

*Proof:* Since $Q_X^{\perp}$ is invariant under rotation, one can express $Q_X^{\perp}$ as the Cholesky factor of $I - Q_X Q_X^{\dagger}$; i.e., $Q_X^{\perp}(Q_X^{\perp})^{\dagger} = I - Q_X Q_X^{\dagger}$. Now, $\Phi Q_X^{\perp}(Q_X^{\perp})^{\dagger} \Phi^{\dagger} = I - (\Phi Q_X)(\Phi Q_X)^{\dagger} = (\Phi Q_X)^{\perp}((\Phi Q_X)^{\perp})^{\dagger}$, and hence the proof. $\square$

This property implies that if the conditions of Proposition 1 are satisfied for the proto-constellation $\mathcal{C}$, then the rotation-based constellation possesses the same antipodal symmetry.

The rotation-based technique is suboptimal in general, because enforcing the associated structure reduces the number of degrees of freedom; $(L-1)T^2 + 2M|\mathcal{C}|(T-M)$ for rotated designs as opposed to $2ML|\mathcal{C}|(T-M)$ for direct designs. Fig. 5 illustrates the distance spectra of 256-point constellations for a system with $M = 2$ and $= 4$ that are generated by the greedy, direct and rotation-based techniques. In order to illustrate the impact of the proto-constellation size on the performance, Fig. 5 shows two rotation-based constellations: the first is generated by 32 rotations of an 8-point proto-constellation, whereas the second is generated by 16 rotations of 16-point proto-constellation. In the greedy and direct cases, the minimum distance of the constellation is 0.7036 and 0.8288, respectively. For the $32 \times 8$ rotation-based design the minimum distance is 0.7633, and that for the $16 \times 16$ design is 0.7113. The decrease in the minimum distance in the $16 \times 16$ case is to be expected, since a larger proto-constellation imposes a more stringent structure on the final constellation. In terms of performance in a system with $N = 2$ receiver antennas, it can be seen from Fig. 5 that at a block error rate of $10^{-5}$ the $16 \times 16$ and $32 \times 8$ rotation-based constellations offer a gain of about 0.2 and 0.4 dB over the greedy constellation, respectively, whereas the directly designed constellation offers a gain of about 0.9 dB. However, the main advantage of rotation-based constellations is that, in comparison with the other two techniques, these constellations

(a) Greedy, $d_{\min} = 0.7036$

(b) Direct, $d_{\min} = 0.8288$

(c) 32 Rotations of 8-pt con-
stell., $d_{\min} = 0.7633$

(d) 16 Rotations of 16-pt con-
stell., $d_{\min} = 0.7113$
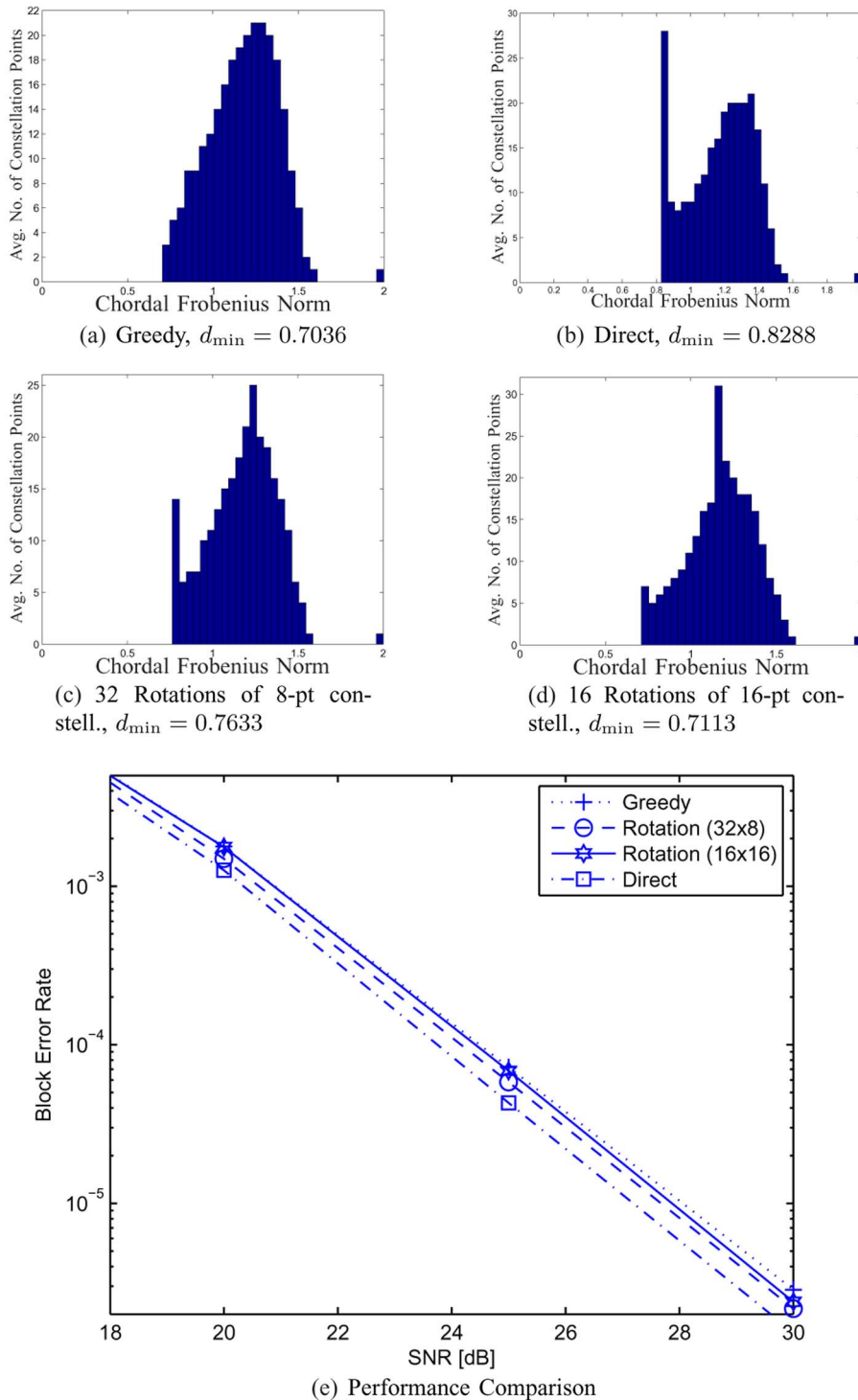
(e) Performance Comparison

Fig. 5. Distance spectra and performance of 256-point constellations for a system with $M = 2$, $T = 4$, and $N = 2$.

are significantly easier to design and store. As we will explain in Section IV-D, the inherent structure of the rotation-based constellations also facilitates the binary labeling of the constellation points. (Although we will not attempt to do so here, we suspect that one might also be able to utilize the structure of rotation-based constellations in order to reduce the detection complexity at the receiver.)

In order to provide some insight into the computational advantage of generating Grassmannian constellations using the ro-

tation-based technique rather than the direct technique, we recall that the accuracy of the current solution of the conjugate-gradient (CG) algorithm with resets on the Grassmann and Stiefel manifolds [21] doubles in at most $J$ iterations, where $J$ is the dimension of the manifold [21, Sec. 3.5.1]. Using this observation, one can use the dimension of the underlying manifolds to compare the number of iterations required to guarantee a certain accuracy in both the direct and the rotation-based designs. In particular, if $\mathcal{C}_L$ is a constellation of the desired cardinality,

$\mathcal{C}$ is the proto-constellation, and $L = |\mathcal{C}_L|/|\mathcal{C}|$ is the number of rotations, then the dimension of the manifold that underlies the direct design is $2ML|\mathcal{C}|(T - M)$, whereas the dimension of the manifold that underlies the rotation-based design is $(L-1)T^2$.[4] As an example, consider the case in which a 256-point constellation is designed using $L = 16$ rotations of a proto-constellation of cardinality $|\mathcal{C}| = 16$ with $T = 2M = 4$. For this case, the manifold that underlies the direct design is of dimension 2048, whereas the manifold that underlies the rotation-based design is of dimension 240.

### D. Quasi-Set-Partitioning Labeling

In order to use Grassmannian constellations in a practical coded communication system, one typically needs to assign a binary label to each point of the constellation. The way in which the points are labeled can have a significant effect on the performance of the system. However, binary labeling of the points in a Grassmannian constellation is difficult, because even for small dimensions, "good" constellations are not known to possess a structure that could be exploited to determine an appropriate labeling strategy. In addition, the number of (real) dimensions of the Grassmann manifold is $2M(T - M)$, which can be quite large for practical signaling scenarios. This large dimensionality renders labeling quite cumbersome, even for, ostensibly plain, Euclidean spaces, let alone Grassmann manifolds. Numerical optimization of the mapping is, in principle, an option (e.g., [29]), but there are $|\mathcal{C}_L|!$ possible labelings, and hence the optimization is a computationally formidable task for all but the smallest constellations. However, we will now show how the inherent structure of rotation-based Grassmannian constellations can be exploited to develop a labeling technique that adheres, to a large extent, to the principles that underlie the standard set-partitioning technique, and hence may be of interest in the development of trellis-coded modulation schemes with Grassmannian signaling; e.g., [30].

We begin by observing that, roughly speaking, standard set-partitioning assigns labels with small Hamming distances to points that lie at large distances in the signaling space. In our case, the signaling space corresponds to a compact Grassmann manifold. Now, if the proto-constellation, $\mathcal{C}$, is properly designed, points in this constellation will lie at maximum pairwise distance. Furthermore, the compactness of the Grassmann manifold implies that introducing more constellation points does not increase the minimum pairwise distance. (It typically reduces it.) From Property 1, we know that rotation preserves the distance between points in the proto-constellation. Hence, the smaller distances in the final constellation, $\mathcal{C}_L$ in (31), occur between points that belong to different rotations of $\mathcal{C}$. Using this insight, we now describe our labeling strategy.

Consider a rotation-based constellation with $|\mathcal{C}| = 2^{n_1}$ and $L = 2^{n_2}$. It is required to label the points of the constellation with binary vectors of length $n_1 + n_2$. For each point in the constellation we will use the first $n_1$ bits to index the point on the underlying proto-constellation, and the remaining bits to index the rotation. By partitioning the label in this way, we ensure that constellation points generated by the same rotation, which will be well-spaced (so long as the proto-constellation is well-designed), differ by a Hamming distance of at most $n_1$ bits. The remaining $n_2$ bits label the rotation, and since there is no known structure for these rotations, these bits can be chosen pseudo-randomly. In general, small proto-constellations provide more degrees of design freedom (in a geometric sense), whereas large proto-constellations endow the final constellation with more structure. Hence, the choice of the cardinality of the proto-constellation and the number of rotations provide a tradeoff between favourable geometric and Hamming distance properties of the constellation.

## V. DETECTION

Having established design principles for the transmitter, we now turn our attention to the receiver. After briefly discussing conventional noncoherent Maximum Likelihood (ML) detection [3], we will describe the proposed reduced search detector. Throughout this section we assume that the channel symbols are drawn from an isotropically distributed Grassmannian constellation.

### A. Maximum Likelihood Detection

For ML detection, we observe from (3) that conditioned on $Q_X$, the received signal $Y$ is a zero-mean isotropically distributed Gaussian random matrix. Hence

$$
\begin{aligned}
&p(Y|Q_X) \\
&= \frac{\exp\left(-\mathrm{Tr}\left(Y^\dagger \left(\frac{M}{\rho T}I_T + Q_X Q_X^\dagger\right)^{-1} Y\right)\right)}{\pi^{TN} \det^N\left(\frac{M}{\rho T}I_T + Q_X Q_X^\dagger\right)} \\
&= \frac{\exp\left(-\frac{\rho T}{M}\mathrm{Tr}\left(Y^\dagger\left(I_T - \frac{1}{1+M/\rho T}Q_X Q_X^\dagger\right)Y\right)\right)}{(\pi M/\rho T)^{TN}(1 + \rho T/M)^{MN}}. \quad (34)
\end{aligned}
$$

A maximum likelihood detector tests the entire constellation, $\mathcal{C}$, in search for the constellation point $\hat{Q}_X$ that maximizes $p(Y|Q_X)$ in (34). Equivalently

$$
\hat{Q}_X = \arg \max_{Q_X \in \mathcal{C}} \mathrm{Tr}\left(Y^\dagger Q_X Q_X^\dagger Y\right). \quad (35)
$$

The main drawback associated with the ML detector is the computational cost of having to examine all possible constellation points in the constellation. In order to increase the computational efficiency, the reduced search algorithm proposed below selects a particular set of candidate points to be examined against the maximum likelihood metric in (35).

### B. Reduced Search Quasi-ML Detection

The reduced search algorithm is based on two concepts: the structure of isotropically distributed Grassmannian constellations and the nature of the received signal. In order to visualize the structure of an isotropically Grassmannian constellation, it is instructive to consider a low-dimensional example. Consider the Grassmann manifold $\mathbb{G}_1(\mathbb{R}^3)$, which is the set of all pairs
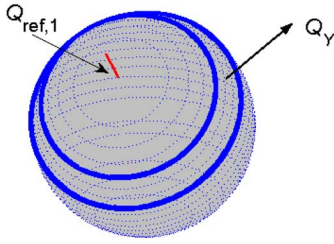
---

[4]In general, the design complexity of the proto-constellation is much less than that of the rotations and hence is ignored in this comparison.

Fig. 6. Reduced Search detection: The width of the band is determined by $A_Y$ and $B_Y$.
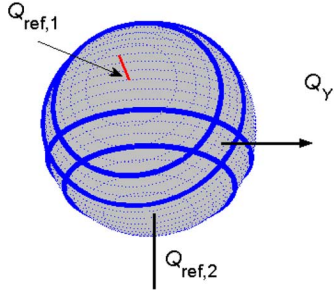


Fig. 7. Pictorial view of the reduced search algorithm when two reference points are used.

of antipodal points that lie on (the surface of) the sphere $S^2$ in $\mathbb{R}^3$. As illustrated in Fig. 6, by observing the manifold from an arbitrary reference point, $Q_{\text{ref},1}$, one can define a sequence of disjoint sets $\mathcal{B}(A_r, A_{r+1}, Q_{\text{ref},1})$ that cover the entire manifold. That is,

$$\mathbb{G}_M(\mathbb{C}^T) = \bigcup_r \mathcal{B}(A_r, A_{r+1}, Q_{\text{ref},1}) \tag{36}$$

where

$$\mathcal{B}(A_r, A_{r+1}, Q_{\text{ref},1}) = \Big\{ Q \in \mathbb{G}_M(\mathbb{C}^T) | A_r$$
$$\leq d_D(Q, Q_{\text{ref},1}) < A_{r+1} \Big\} \tag{37}$$

and $d_D(\cdot, \cdot)$ is a distance metric (not necessarily the chordal Frobenius norm used in Section IV) and $\{A_r\}$ is a set of appropriate threshold values. As illustrated in Fig. 6, each of these sets constitutes a 'band' on the Grassmannian manifold. The partitioning of the manifold in (36) suggests that the points in a Grassmannian constellation can be classified according to the band in which they lie, and we will exploit that partitioning in the reduced-search detector proposed below.

For a perfect Grassmannian constellation, the elements of the set $\{A_{r+1} - A_r\}$ can be chosen sufficiently small for each band $\mathcal{B}(A_r, A_{r+1}, Q_{\text{ref},1})$ to either contain no constellation points or to contain points that lie on a "circle" of a certain radius. For a practical constellation that is not perfectly uniform, the corresponding constellation points might be slightly perturbed and might not necessarily lie on the contour of a circle, but they will lie within a specific band.

In addition to the Grassmannian structure, the reduced search algorithm is also based on a specific feature of the received signal, $Y$. In particular, we have the following result.

*Theorem 2:* Let $Y = Q_Y R_Y = Q_X H + V$, where $Q_Y R_Y$ is a nonunique QR decomposition of the received signal $Y$, $Q_X$ is a $T \times M$ isotropically distributed random unitary matrix, $V \in \mathbb{C}^{T \times N}$ is an isotropically distributed random Gaussian matrix and $H \in \mathbb{C}^{M \times N}$. Then the mutual information between $Q_X$ and $R_Y$ is zero and the mutual information between $Q_Y$ and $H$ is zero. That is,

$$I(Q_X; R_Y) = 0 \tag{38}$$

and

$$I(H; Q_Y) = 0. \tag{39}$$

*Proof:* See Appendix E. □

Theorem 2 states that the perturbation in the signal subspace caused by the isotropically distributed additive noise does not couple the information about the signal subspace and the channel state. That is, all the information about the subspace spanned by the columns of $Q_X$ is contained in the subspace spanned by the columns of $Q_Y$ and all the channel state information is contained in $R_Y$.

In order to efficiently generate reliable decisions, the reduced search algorithm decomposes the ML detection in (35) into two steps. In the first step the detector uses the information contained in $Y = Q_Y R_Y$ and the look-up table to select a set of candidate constellation points. In the second step the detector selects the constellation point in that set with the largest likelihood; cf. (35). Based on Theorem 2, the information contained in $R_Y$ and the information contained in $Q_Y$ can be processed separately in the first step. Indeed, in the proposed reduced search algorithm, for each received signal matrix $Y$, the information contained in $R_Y$ is used to determine the size of the search region and the information contained in $Q_Y$ to determine to the location of the search region.

We now describe the detection procedure in more detail. The reduced search detector employs a look-up table that contains the distance between all the constellation points $\{Q_{X_i}\}_{i=1}^{|\mathcal{C}|}$ and the reference point $Q_{\text{ref},1}$. This look-up table is only constructed once, prior to implementation, and is stored at the receiver. At each channel use, the receiver initiates the detection process by computing the QR decomposition of $Y$. Since Theorem 2 states that all the information about $Q_X$ is contained in $Q_Y$, the receiver attempts to localize the search by computing the distance between $Q_Y$ and the reference point, $d_D(Q_Y, Q_{\text{ref},1})$. Using this distance, the receiver then consults the look-up table to select a certain set of candidate constellation points. To determine the size of this set, the receiver uses the channel information contained in $R_Y$ to generate two real values, $A_Y$ and $B_Y$. Observing the manifold from $Q_{\text{ref},1}$, these values are used to define a band that contains $Q_Y$, $\mathcal{B}(d_D(Q_Y, Q_{\text{ref},1}) + A_Y, d_D(Q_Y, Q_{\text{ref},1}) + B_Y, Q_{\text{ref},1})$. The reduced set of candidate constellation points is then defined as the set of points in this band. That is, if we denote this set by $\mathcal{C}'(A_Y, B_Y, Q_{\text{ref},1})$, we have

$$\mathcal{C}'(A_Y, B_Y, Q_{\text{ref},1}) = \{Q_X \in \mathcal{C} | A_Y \leq d_D(Q_X, Q_{\text{ref},1})$$
$$- d_D(Q_Y, Q_{\text{ref},1}) < B_Y\}. \tag{40}$$

Having determined $\mathcal{C}'(A_Y, B_Y, Q_{\mathrm{ref},1})$, the second step of the detection process involves the examination of all the constellation points that lie within this set against the maximum likelihood metric in (35). That is, the receiver decides in favour of

$$\check{Q}_X = \arg \max_{Q_X \in \mathcal{C}'} \mathrm{Tr}\left(Y^\dagger Q_X Q_X^\dagger Y\right). \qquad (41)$$

Since $\mathcal{C}'$ is a subset of $\mathcal{C}$, and since the cardinality of $\mathcal{C}'$ depends on the width of the band in (40), if the width of the band is decreased, the number of likelihood computations is reduced but the probability of missing the correct constellation point is increased. On the other hand, if the width of the band is increased, the probability of missing the correct constellation point is reduced at the expense of computing redundant likelihoods. The tradeoff between the cardinality of $\mathcal{C}'$ and performance can be controlled through the choice of the threshold values $A_Y$ and $B_Y$. In Section VII we will show that by carefully choosing $A_Y$ and $B_Y$, the cardinality of $\mathcal{C}'$ can be considerably reduced from that of $\mathcal{C}$ without significant loss in performance.

A further step in applying the reduced search detection strategy is to augment the look-up table by including distances from the constellation points to several other reference points. In particular, if distances from another reference point $Q_{\mathrm{ref},2}$ are recorded in the look-up table, one can measure $d_D(Q_Y, Q_{\mathrm{ref},1})$ and $d_D(Q_Y, Q_{\mathrm{ref},2})$ and only consider the candidate points $Q_{\hat{X}_k}$ such that

$$Q_{\hat{X}_k} \in \bigcap_{i=1,2} \mathcal{C}'(A_Y, B_Y, Q_{\mathrm{ref},i}).$$

From Fig. 7, it is evident that the the volume of the intersection of the two bands can be significantly less than the volume of one band. However, this reduction in the search space must be properly accounted for in the computation of the values of $A_Y$ and $B_Y$. In Section VII, the potential impact of using several reference points on reducing the search space is investigated.

One appropriate choice for the values of $A_Y$ and $B_Y$ in (40) would be those that yield the smallest value of the band width $|A_Y - B_Y|$ that guarantees that the probability that the correct constellation point does not lie inside $\mathcal{C}'$ is less than a small number, say $\delta \geq 0$. That is,

$$(A_Y, B_Y) = \arg \inf_{\{(A,B)|P(Q_X \notin \mathcal{C}'|Y) < \delta\}} |A - B| \qquad (42)$$

where $\mathcal{C}'$ is defined in (40). Direct computation of the probability $P(Q_X \notin \mathcal{C}'|Y)$ is quite complicated and depends on the choice of the reference point. As an alternative we propose to use Chebychev's inequality to bound this probability. Doing so, in Appendix F we obtain closed form expressions for $A_Y$ and $B_Y$ that can be used to generate some insight into the characteristics of the reduced search algorithm. In agreement with Theorem 2, we will show that these threshold values are independent of both the transmitted signal subspace spanned by the columns of $Q_X$ and the received signal subspace spanned by the columns of $Q_Y$. That is, the values of $A_Y$ and $B_Y$ depend only on $R_Y$. In fact, these values are dominated by the smaller singular values of the matrix $R_Y$. (In Appendix F, we will show that $R_Y$ is closely related to the channel matrix $H$.) By adjusting the values of $A_Y$

and $B_Y$ one can ensure that the correct constellation point belongs to the set of candidate constellation points with some prescribed high probability. For a channel matrix with small singular values, the width of the band, $|A_Y - B_Y|$, is typically large, which indicates that exhaustive search maximum likelihood detection is required, whereas for channel matrices with large singular values, the width of the band is small, which indicates that only a few constellation points need to be tested against the ML metric.

In Appendix F, we have shown that if the probability that the correct constellation point lies outside the search region is held constant, the width of the search region decays at least as fast as $1/\sqrt{\rho}$. However, in Appendix F we have argued that it is desirable for this probability of missing the correct constellation point to decay with the SNR. In Appendix F we propose a particular way for choosing the parameters of the reduced search algorithm so that this probability decays with the SNR. For these choices, the width of the search region decay s as $\log(\rho)/\sqrt{\rho}$, irrespective of $T$, $M$ and $N$, and also irrespective of the number of reference points and the cardinality of the constellation.

One drawback of the reduced search scheme is the need to store the look-up table increases the memory requirement of the receiver. This is especially true if many reference points are used. This drawback can be reduced by quantizing distances from the reference points to the constellation points. For example, the look-up table can be organized such that the $ij$th cell holds the index of the constellation points that lie at distances between $d_i$ and $d_{i+1}$ from the $j$th reference point. Using this approach we can obtain valuable reduction in the required memory.

## VI. PAIRWISE ERROR PROBABILITY

In the previous sections, we considered the design and detection of Grassmannian constellations that enable the high SNR ergodic capacity of the MIMO system to be approached. In this section, we consider the performance of such constellations if the data rate is fixed and the SNR is allowed to increase. Since computing an exact expression for the probability of block error is usually infeasible in space–time signaling contexts, the pairwise error probability has often been employed to provide useful insight into the key features that govern the high SNR performance of a given space–time coding scheme [4], [31], [32]. For unitary signaling in communication scenarios in which the channel is not known at the receiver, exact expressions for the pairwise error probability were derived in [4] and [32]. However, the evaluation of these exact expressions is numerically unstable because they involve the computation of residues at poles of high multiplicities [32]. In order to avoid this numerical inconvenience, bounds on the pairwise error probabilities were developed in [4], [32] and [33]. In this section we provide an alternative expression for the exact pairwise error probability. Our method differs from the one presented in [4] and [32] in that it produces an expression in the form of a series expansion and avoids the computation of residues. Our series expansion is absolutely convergent and hence can be used to compute the pairwise error probability up to the required degree of precision. In fact our method will not only yield an expression for the pairwise error probability, it will also give an expression for the

distribution of the multivariate random variable involved in the computation of this probability (see Appendix G). This distribution may be useful in assessing the performance of "soft-decision" based detectors.

The pairwise error probability (PEP) $P(i \to j)$ is defined as the probability that the receiver mistakes the $i$th constellation point, $Q_{X_i}$, for the $j$th constellation point, $Q_{X_j}$, given that the $i$th constellation point has been transmitted. In order to compute $P(i \to j)$, we will define the matrix $\Gamma_{ij}$ to be

$$\Gamma_{ij} = \begin{bmatrix} \gamma_{11} \left(\Sigma_{ij}^2 - I\right) & \gamma_{12} \left(I - \Sigma_{ij}^2\right)^{1/2} \Sigma_{ij} \\ \gamma_{12} \left(I - \Sigma_{ij}^2\right)^{1/2} \Sigma_{ij} & \gamma_{22} \left(I - \Sigma_{ij}^2\right) \end{bmatrix} \quad (43)$$

where $\gamma_{11} = \frac{1}{2}\left(1 + \frac{M}{\rho T}\right)$, $\gamma_{22} = \frac{1}{2}\frac{M}{\rho T}$, $\gamma_{12} = \sqrt{\gamma_{11}\gamma_{22}}$, and $V_{ij}\Sigma_{ij}U_{ij}^\dagger$ is the SVD of $Q_{X_j}^\dagger Q_{X_i}$. Let $\lambda_k^{ij}$ be the $k$th eigenvalue of $\Gamma_{ij}$ and $\kappa(k)$ be the multiplicity associated with this eigenvalue. Let $X_{ij}$ be the random variable defined as

$$X_{ij} = \sum_{k=1}^{K_{ij}} \lambda_k^{ij} \chi_{2\kappa(k)N}^2 \quad (44)$$

where $\chi_n^2$ denotes a Chi-Square random variable with $n$ degrees of freedom and $K_{ij}$ is the number of distinct eigenvalues of $\Gamma_{ij}$. We show in Appendix G that the pairwise error probability $P(i \to j)$ is given by

$$\begin{aligned} P(i \to j) &= P(X_{ij} \geq 0) \\ &= P(X_{ij} \geq x)\big|_{x=0}. \end{aligned} \quad (45)$$

In order to compute $P(X_{ij} \geq x)$, we invoke the following result from [34], [35] to find the distribution function of $X_{ij}$.

*Lemma 6:* Let $X = a(\chi_{2n}^2 + a_1\chi_{2n_1}^2 + \cdots + a_r\chi_{2n_r}^2 - b_1\chi_{2m_1}^2 - \cdots - b_s\chi_{2m_s}^2)$, where the Chi-square variates are independent and $a, a_1, \ldots, a_r, b_1, \ldots, b_s$ are positive constants such that $a_i \geq 1$, $b_i \geq 1$. Define constants $d_j$ and $d_k'$ by the identities,

$$\prod_{i=1}^{r} a_i^{-n_i} \left(1 - \left(1 - a_i^{-1}\right) z\right)^{-n_i} \equiv \sum d_j z^j, \quad (46a)$$

$$\prod_{i=1}^{s} b_i^{-m_i} \left(1 - \left(1 - b_i^{-1}\right) z\right)^{-m_i} \equiv \sum d_k' z^k. \quad (46b)$$

Let $N_f = 2(n + n_1 + \cdots + n_r)$, $M_f = 2(m_1 + \cdots + m_s)$. Then for every $x$,

$$P(X \geq x) = \sum_{j=0}^{\infty} \sum_{k=0}^{\infty} d_j d_k' H_{M_f + 2k, N_f + 2j}\left(\frac{x}{a}\right) \quad (47)$$

where $H_{2n,2m}(x) = \int_{-\infty}^{x} h_{2n,2m}(u)du$

$$h_{2n,2m}(u) = \begin{cases} \sum_{s=0}^{n-1} 2^{-(s+m)} \frac{(m)_s}{s!} f_{2n-2s}(u), & x \geq 0 \\ \sum_{s=0}^{m-1} 2^{-(s+n)} \frac{(n)_s}{s!} f_{2m-2s}(-u), & x < 0 \end{cases}$$

$f_{2n}(x) = \frac{1}{2^n(n-1)!} e^{-x/2} x^{n-1}$, $(p)_n = p(p+1)\ldots(p+n-1)$ and $(p)_0 = 1$. $\qquad\qquad\square$

Using (45) and (47), one can analytically compute

$$\begin{aligned} P(i \to j) &= P(X_{ij} \geq 0) \\ &= \sum_{\ell=0}^{\infty} \sum_{k=0}^{\infty} d_\ell d_k' \\ &\quad \times \sum_{s=0}^{N_f/2+\ell-1} 2^{-(s+M_f/2+k)} \frac{(M_f/2+k)_s}{s!}. \end{aligned} \quad (48)$$

The union bound is therefore given by

$$P_U = \frac{2}{|\mathcal{C}|} \sum_{i=1}^{|\mathcal{C}|-1} \sum_{j=i+1}^{|\mathcal{C}|} P(i \to j). \quad (49)$$

Notice that $d_\ell \geq 0$, $d_k' \geq 0$, $\sum d_\ell = 1$ and $\sum d_k' = 1$ and hence the series in (46) is absolutely convergent. Moreover, for given $T$, $M$ and $N$, the inner summation is independent of both the constellation size and the SNR. Hence the inner summation need only to be computed once and can be used to assess the PEP performance of the communication system at different rates and SNRs. For $T = 2M$ and constellations generated by (24), one can use the results of Proposition 2 in Appendix D, to show that only half the terms of (49) need to be computed.

The PEP expression in (48) can help to identify the impact of changes in the design parameters on the asymptotic performance of the communication system. For instance, the expressions in (44) and (48) suggest that the number of receive antennas, $N$, has an impact on the pairwise error probability only through the number of degrees of freedom of the associated Chi-square random variables. An increase in the degree of freedom results in a Chi-square distribution that is more concentrated around its mean. Hence, if $N$ is increased the areas of overlap between the positively weighted and the negatively weighted Chi-square variables in (44) is reduced, which results in a significantly smaller PEP. Notice that the number of receive antennas does not affect the magnitude nor the multiplicity of the eigenvalues of $\Gamma_{ij}$ in (43).

## VII. NUMERICAL RESULTS

In this section we provide a few numerical results that illustrate the efficacy of our approaches to constellation design and detection. The channel is assumed to be independent Rayleigh block-fading with coherence time $T = 4$. In order to achieve the maximum number of degrees of freedom of this channel (cf. (1) and (2)), the numbers of transmit and receive antennas were chosen to be $M = N = T/2 = 2$. In Section VII-A we will evaluate the performance and the computational cost of the reduced search detector introduced in Section V-B, in Section VII-B we will compare the performance of constellations designed using the techniques described in Section IV against that of some existing unitary constellations, and in Section VII-C we will provide performance comparisons with some training-based signaling schemes. In these simulation experiments, we will consider Grassmannian constellations of size 16, 256, 512, 1024, and 4096, which correspond to
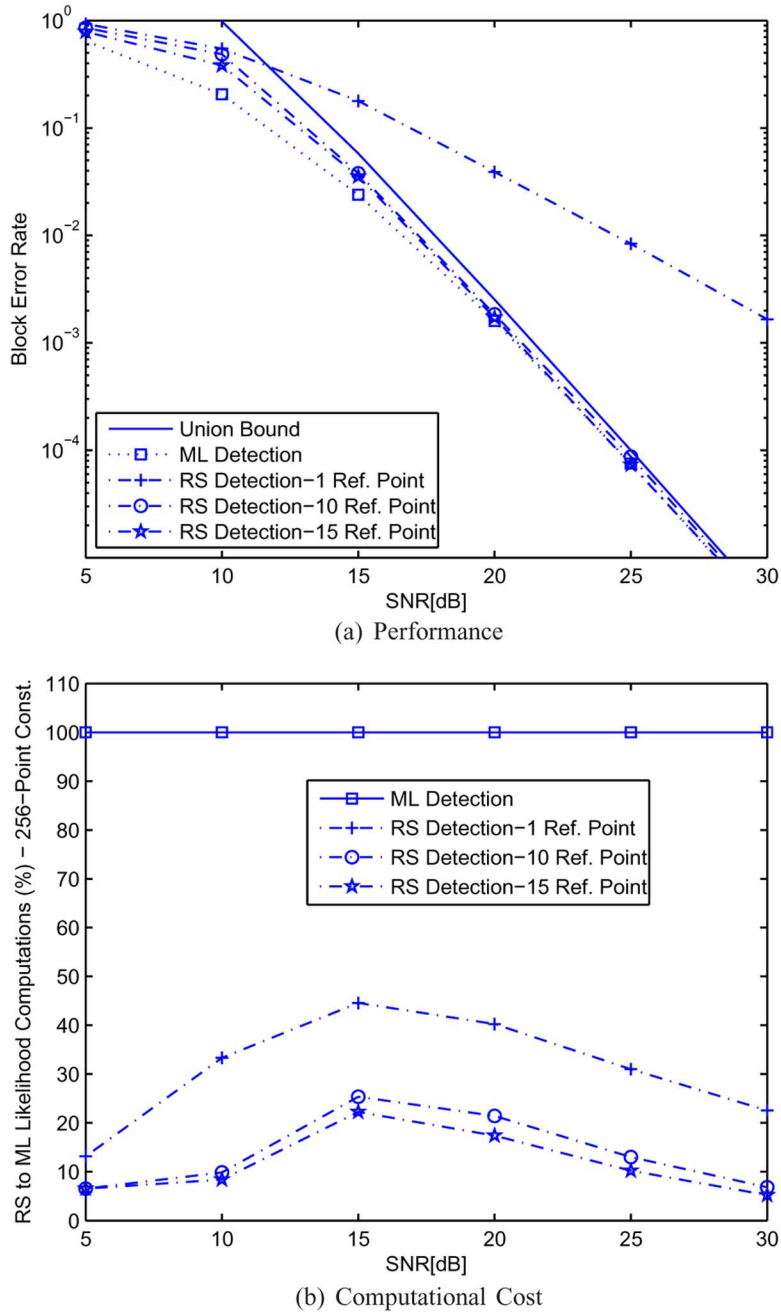
(a) Performance



(b) Computational Cost

Fig. 8. The performance and computational advantage of the Reduced Search detection algorithm of Section V over that of ML detection for the 256-point Grassmannian constellation.

data rates of 1, 2, 2.25, 2.5, and 3 bits per channel use (bpcu), respectively.

### A. Evaluation of the Reduced Search Detector

In order to assess the effectiveness of the reduced search algorithm introduced in Section V, in Figs. 8(a) and 9(a) we have plotted the corresponding block error rate when different numbers of reference points are used, for systems employing the greedily designed constellations of size 256 and 1024. These systems operate at data rates of 2 and 2.5 bpcu, respectively. For all simulations in these figures, the parameter $k$ of reduced search algorithm was chosen according to (81) in Appendix F

with $k_1 = 1$ and $c = 0.25$. In these figures we have also plotted the block error rate of the ML detector, as well as the union bound derived from the exact expression for the pairwise error probability of Section VI. From these figures, we observe that when the reduced search algorithm uses multiple reference points its performance can indeed approach that of the ML detector. In addition, our expression for the pairwise error probability seems to lead to a tighter union bound than the asymptotic union bound derived in [20].

In Figs. 8(b) and 9(b), we have plotted the number of likelihood computations required for the reduced search algorithm. From these figures, one can conclude that an increase in the

Fig. 9. The performance and computational advantage of the reduced search detection algorithm of Section V over that of ML detection for the 1024-point Grassmannian constellation.
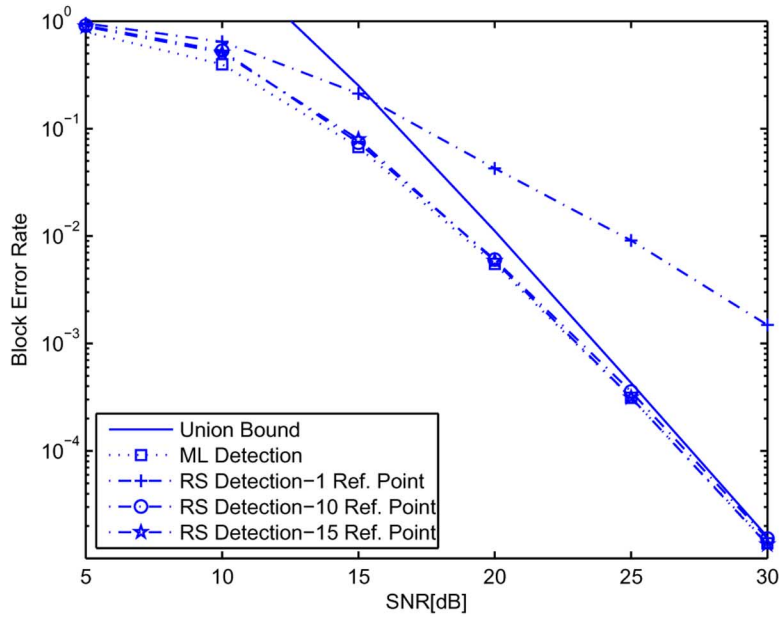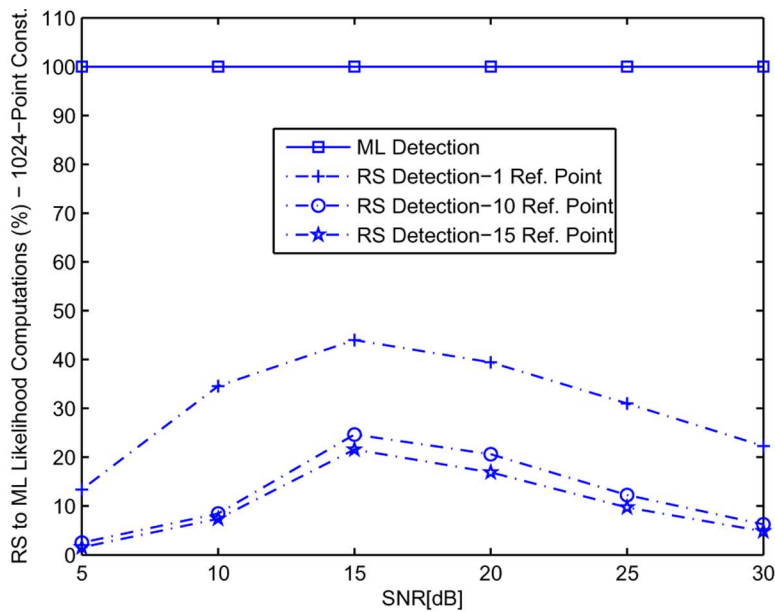
number of reference points not only results in the performance of the reduced search algorithm approaching that of the ML detector; there is also a significant reduction in the number of likelihood computations. For example, at an SNR of 30 dB, when the detector uses 15 reference points instead of 1 reference point, the average number of likelihood computations is reduced from about 56.3 to 15.4 evaluations for the 256-point constellation and from about 225.3 to 51.2 likelihood computations for the 1024-point constellation. (ML detection requires the evaluation of 256 and 1024 likelihoods, respectively.) Observe that because we choose $k$ as in (81), the width of the search region (and consequently the number of candidates to be considered) initially

expands with the SNR, and then decays. The asymptotic decay of the width takes the form $\log(\rho)/\sqrt{\rho}$; see Appendix F for further discussion.

### B. Performance Comparison of Greedily Designed Constellations

In this section we provide performance comparisons, based on (full) ML detection, between the Grassmannian constellations designed using the greedy technique in Section IV-A and those designed in [20], which we will refer to as the MBV constellations. The design of the MBV constellations is based on greedily minimizing the asymptotic union bound
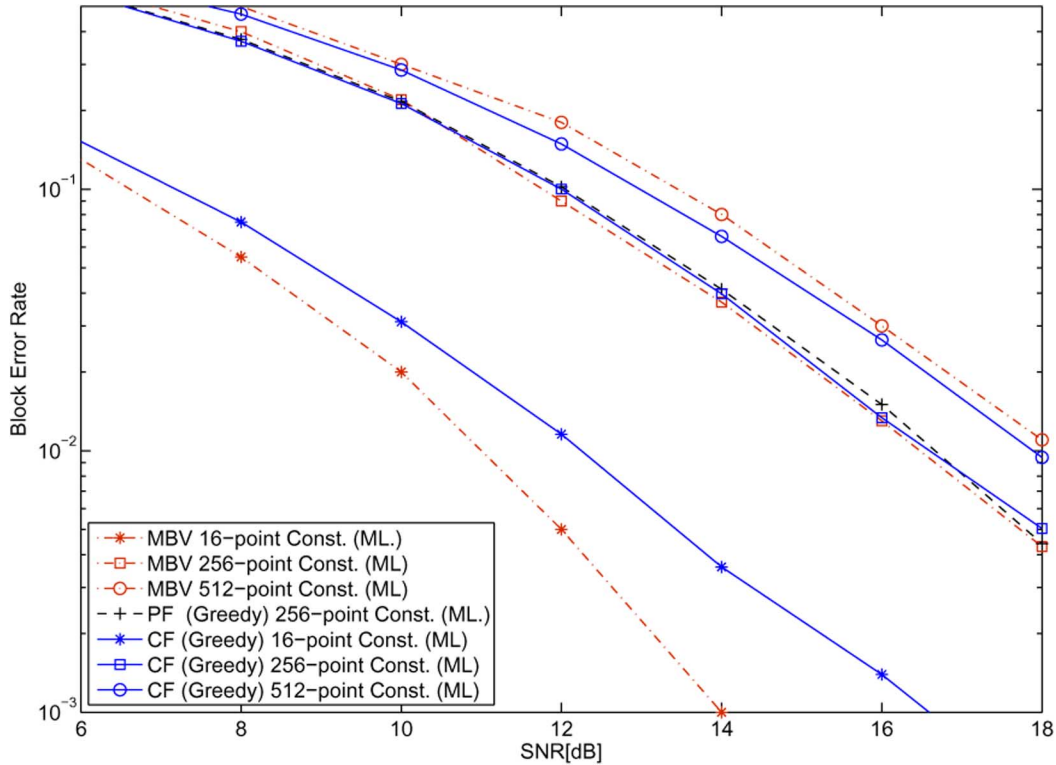
Fig. 10. Performance comparison between constellations designs via greedy techniques that use the chordal Frobenius (CF) norm, the asymptotic union bound (MBV) and the projection Frobenius (PF) norm as the distance metric. ML detection is used for all constellations.

derived in [20], whereas our approach to constellation design (cf. Section IV) is based on maximizing the chordal Frobenius norm between constellation points. The performance results in Fig. 10 show that for 16-point constellations (i.e., a data rate of 1 bpcu) the MBV constellation outperforms our greedily designed constellation by about 2.5 dB at a block error probability of $10^{-3}$. When the data rate is increased to 2 bpcu, the 256-point MBV constellation provide better performance than our 256-point constellation, but only slightly better. When the data rate is further increased to 2.25 bpcu (i.e., for the case of 512-point constellations), our (greedy) constellation performs better than the MBV constellations. This figure shows that when the system is operating at higher data rates our greedily designed constellations can outperform the (greedy) MBV ones.

Fig. 10 also shows the performance of a 256-point constellation that was generated using the greedy algorithm, but with the projection Frobenius norm in (21) instead of the chordal Frobenius norm proposed herein; cf. (20). It can be seen from this figure, that the constellation designed using the chordal Frobenius norm performs slightly better than that designed using the projection Frobenius norm for SNRs up to about 17 dB. However, the constellation designed using the projection Frobenius norm performs slightly better at higher SNRs. This suggests that for systems operating close to the ergodic capacity, the chordal Frobenius norm is a more appropriate distance metric than the projection Frobenius norm. For systems operating away from the capacity limit, a distance metric that directly accounts for the probability of error may be more appropriate.

Finally, we emphasize that the comparisons in this section have focused on the performance of greedily designed constellations. Better performance can be achieved by constellations that are designed using the joint techniques outlined in Sections IV-B and C; cf. Figs. 4 and 5.

### C. Comparison With Training-Based Schemes

In [7] a training-based signaling scheme for the noncoherent MIMO channel was presented. In the training phase of this scheme, the transmitter uses the channel $M$ times to send pilot symbols. The receiver uses these pilot symbols to generate a minimum mean square (MMSE) estimate of the channel. Assuming this estimate to be sufficiently accurate, the receiver then coherently detects the data sent by the transmitter during the remaining $T - M$ channel uses, using a 'mismatched' ML detector. This phase is known as the data communication phase.

As a first comparison of the performance of the training-based scheme with that of Grassmannian signaling, we consider systems operating at rates of 1, 2, and 2.5 bpcu. For Grassmannian signaling these rates correspond to constellations of size 16, 256, and 1024, respectively. For this comparison, we have used constellations that were designed using the direct technique in Section IV-B. For the training-based scheme we have used Alamouti signaling [36] in the data communication phase. To match the considered rates, the underlying constellations of the training-based scheme were chosen to be 4, 16, and 32-QAM, respectively. (For Alamouti signaling, 'mismatched' ML detection can be achieved by linear processing and (low-complexity) symbol-by-symbol detection.) From Fig. 11, it can be seen
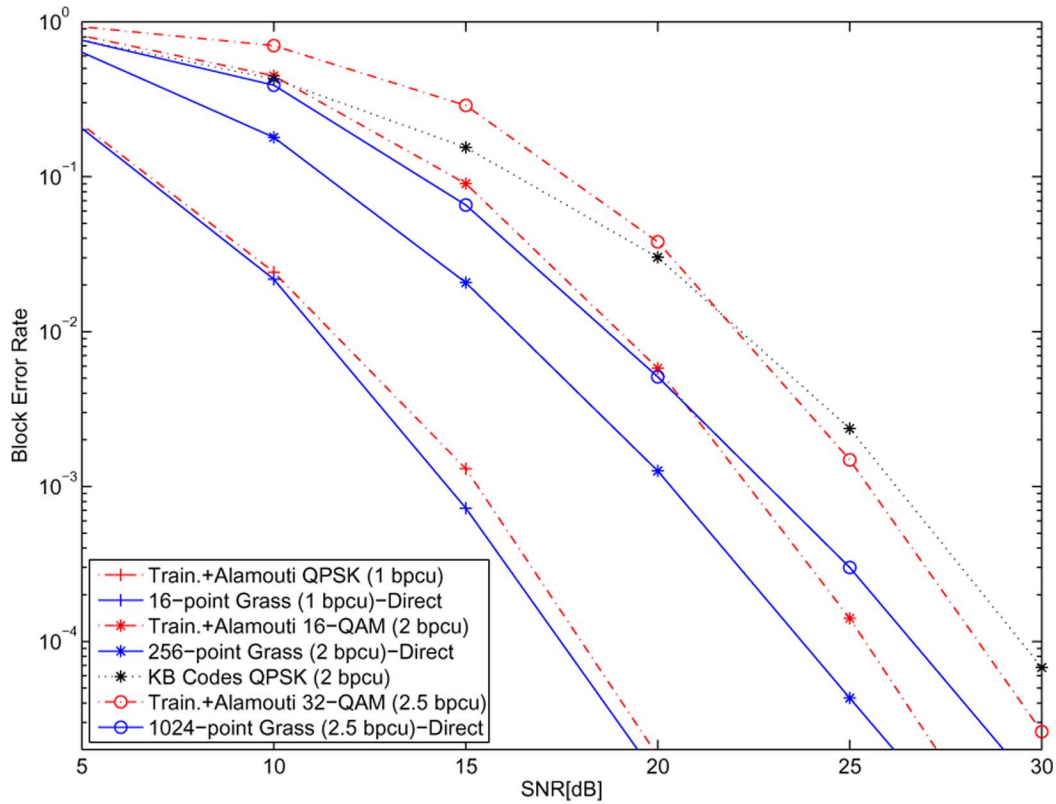
Fig. 11. Performance comparison between signaling using the proposed constellations with ML detection and Alamouti-based training signaling schemes with "mismatched" detection. The Grassmannian constellations in this figure were designed using the direct method of Section IV-B. The performance of the 256-point constellation in [18] (denoted by KB) with ML detection is also shown.

that the Grassmannian constellations that are designed using the direct technique of Section IV-B perform better than the Alamouti-based training schemes. Furthermore, Fig. 11 shows that the advantage of signaling with Grassmannian symbols over signaling with Alamouti-based training increases with the data rate. This advantage directly follows from the asymptotic optimality of Grassmannian signaling at high SNR.

In Fig. 11 we have also plotted the performance of the 256-point constellation proposed in [18], which is based on an exponential mapping of a coherent space–time constellation. One can see from this figure, that due to the stringent structure imposed by the design methodology in [18], both the Alamouti-based training scheme and the proposed Grassmannian signaling scheme provide better performance.

Although the Alamouti scheme is attractive from a detection complexity perspective, it is known to be rate-suboptimal when the number of receive antennas, $N$, is greater than one [37]. In the $N = 2$ case that we are considering, a constellation that exhibits more favourable performance characteristics is that in [38], which we will refer to as the Golden constellation. In Fig. 12 we compare the performance of a 4096-point Grassmannian constellation (designed using 512 rotations of an 8-point directly designed proto-constellation) with that of a training-based scheme that utilizes a Golden constellation with 8-QAM symbols. (Both schemes have a data rate of 3 bpcu.) In this figure we use ML detection for the Grassmannian constellations, and both 'mismatched' and 'optimal' detection [39] for the detection of the Golden constellation. At a block

error rate of $10^{-2}$, one can see that the proposed constellation has an SNR advantage of about 1.5 dB over the training-based scheme with the Golden constellation and 'optimal' detection, and about 2 dB over the same scheme when 'mismatched' detection is used. While the complexity of mismatched detection of Alamouti-based training schemes is quite low, detection of Golden constellations is significantly more expensive. In particular, the number of multiplications required for 'optimal' and 'mismatched' detection of the Golden constellation is $O(TNM|\mathcal{C}|) + O(NM^2|\mathcal{C}|)$ and $O(T_d NM|\mathcal{C}|)$, respectively, where $T_d$ is the time interval of the data communication phase. For comparison, the complexity of ML detection of noncoherent Grassmannian symbols is $O(N(T^2 + T)|\mathcal{C}|)$. Fig. 12 also shows the performance of a randomly generated Grassmannian constellation. The points of this constellation are the unitary components of the QR decomposition of $T \times M$ matrices with independent and identically distributed complex zero mean Gaussian entries. From Fig. 12 it can be seen that although our 'carefully' designed constellations perform better than the randomly chosen constellation, the difference is not very large. However, our rotation-based constellation has the additional advantage that it possesses a structure that reduces its storage requirements and enables a quasi-set-partitioning labeling scheme; cf. Section IV-D.

In order to develop further insight into the performance of the proposed constellation designs, we recall that it was shown in [2] that at high SNRs, every 3 dB increase in the SNR results in a capacity gain of $M(1 - M/T) = 1$ bpcu. Now, by com-
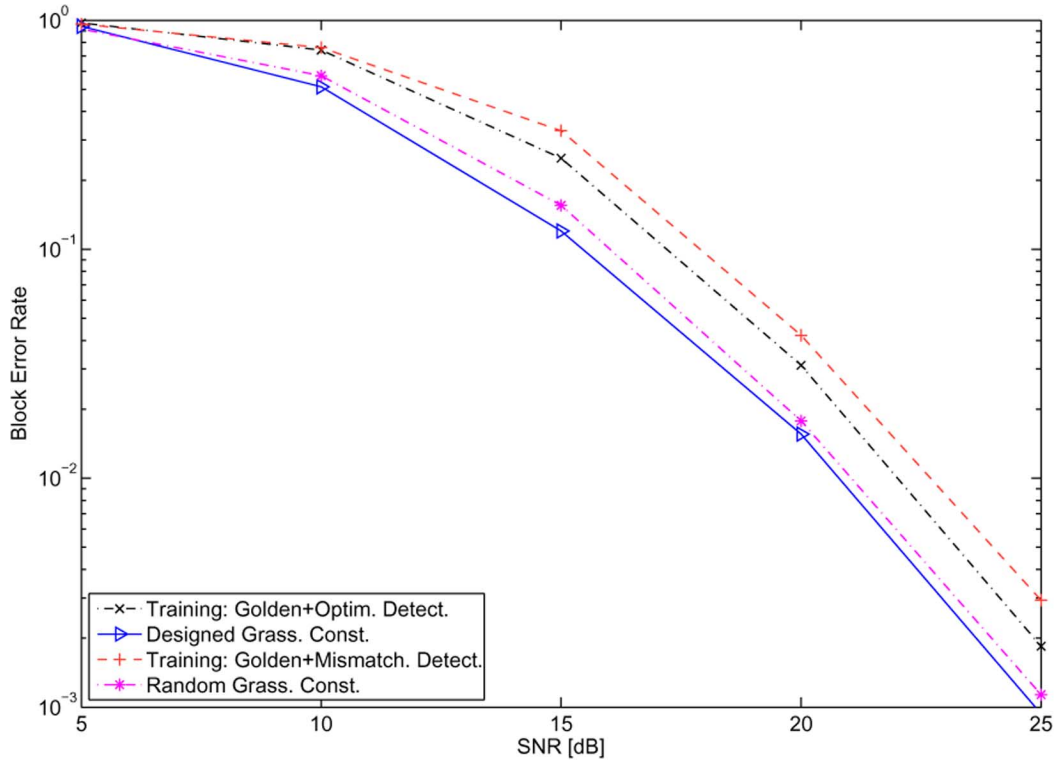
Fig. 12. Comparison between the performance of the proposed codes with ML detection and a training-based scheme with Golden code and "mismatched" and "optimal" detection. The performance of randomly generated Grassmannian constellations with ML detection is also shown.

paring the performance of our 256-point constellation (2 bpcu) in Fig. 11 with that of our 4096-point constellation (3 bpcu) in Fig. 12 at a fixed block error rate, say $10^{-2}$, one can see that the SNR gap between these constellations is only slightly more than 3 dB, which is the difference that one would predict from the capacity expression in [2]. This suggests that our constellations will enable the high SNR ergodic noncoherent capacity of the MIMO system to be approached, and is due, in part, to the fact that these constellations were designed using a metric that appropriately accounts for the manner in which noise perturbs the signal subspace.

## VIII. CONCLUSION

In this paper, we considered noncoherent MIMO communication over a block Rayleigh fading channel using Grassmannian constellations. We began by studying the manner in which the channel and noise matrices perturb the subspace spanned by the transmitted signal matrix. This perturbation analysis was used to determine an appropriate distance metric for constellation design and to develop an efficient suboptimum detection strategy. The rationale behind choosing the distance metric was to quantify the perturbation that the signal subspace undergoes, in a sense that conforms with the communication model. This metric was found to be the chordal Frobenius norm, rather than the commonly used projection Frobenius norm. Having determined the appropriate distance metric, we then developed greedy, direct and rotation-based techniques for designing Grassmannian constellations. These techniques provide a tradeoff between the

quality of the distance spectrum and design complexity. In addition, the rotation-based technique yields constellations that are easy to store and regenerate, and admit a binary labeling of the constellations that resembles, to a large extent, standard set-partitioning labeling. As we illustrated in our simulations, the choice of the appropriate metric for constellation design leads to an appreciable performance gain over constellations that are designed according to other distance metrics. Furthermore, the proposed constellations exhibit performance characteristics that conform to established theoretical results.

In addition to guiding the development of the proposed constellation design techniques, the subspace perturbation analysis also enabled us to show that the information about the transmitted constellation point is captured by the $Q$-component of the received signal matrix whereas all the channel information is contained in the $R$-component. We then used this result to introduce an efficient suboptimum detection strategy with a search region whose width decays as $\log(\rho)/\sqrt{\rho}$, where $\rho$ is the SNR. Finally, we derived an exact expression for the pairwise error probability. This expression is numerically stable and was shown to be significantly tighter than other approximate expressions that have been developed in the literature.

The analysis herein of the detection strategy and the pairwise error probability could prove valuable in identifying the role of various design parameters. For instance, while the capacity gain of choosing $N > M$ was discussed in [2], using the analyses in this paper we are able to assess the potential gains of choosing $N > M$ from other perspectives. In particular, we have shown in Section VI how the increase in the number of receive antennas

is reflected in the asymptotic error performance of the communication system. Similarly, in Appendix F-B we have shown that this increase can offer an attractive reduction in the computational complexity of the proposed detector.

## APPENDIX A
## PROOF OF LEMMA 2

Let $U_X = [Q_X \quad Q_X^\perp]$. In order to show the desired result, we will use the following lemma.

*Lemma 7:* If $G$ and $\hat{G}$ are defined as in (9), then $G$ and $\hat{G}$ are independent of each other and independent of $U_X$.

*Proof:* Let $\Theta = [G^\dagger \quad \hat{G}^\dagger]^\dagger$. Since the additive noise matrix $V$ in (3) is isotropically distributed then $\tilde{V} = U_X^\dagger V = \Theta$ is also isotropically distributed with independent complex Gaussian entries. This implies that $G$ and $\hat{G}$ are independent of each other. Furthermore, we have that

$$
\begin{aligned}
p(\Theta|U_X) &= p\left(U_X^\dagger V|U_X\right) \\
&= \alpha \exp\left(-\operatorname{Tr}\left(V^\dagger U_X U_X^\dagger V\right)\right) \\
&= \alpha \exp\left(-\operatorname{Tr}(V^\dagger V)\right) = p(V)
\end{aligned}
$$

and hence $G$ and $\hat{G}$ are independent of $U_X$.  □

Using Lemma 7, the entries of

$$
\left(H + \sqrt{\frac{M}{\rho T}} Q_X^\dagger V\right) \sim \mathcal{CN}\left(0, 1 + \frac{M}{\rho T}\right)
$$

and

$$
\sqrt{\frac{M}{\rho T}}\left(Q_{X_1}^\perp\right)^\dagger V \sim \mathcal{CN}\left(0, \frac{M}{\rho T}\right)
$$

are statistically independent. If we let

$$
B = \left[\left(H + \sqrt{\frac{M}{\rho T}} G\right)^\dagger \quad \sqrt{\frac{M}{\rho T}} \hat{G}_1^\dagger\right]^\dagger
$$

as in (11), we can write

$$
B = \Lambda \tilde{B} \tag{50}
$$

where $\Lambda = \left[\sqrt{1 + \frac{M}{\rho T}} I_M \quad \sqrt{\frac{M}{\rho T}} I_{N-M}\right]$ and $\tilde{B}$ is a random matrix with i.i.d. Gaussian entries. Now, for $N > M$ we have that

$$
\begin{aligned}
\frac{1}{\lambda_{\min}(BB^\dagger)} &= \lambda_{\max}\left((BB^\dagger)^{-1}\right) \\
&= \lambda_{\max}\left(\Lambda^{-2}(\tilde{B}\tilde{B}^\dagger)^{-1}\right) \\
&\leq \frac{1}{\frac{M}{\rho T}\lambda_{\min}(\tilde{B}\tilde{B}^\dagger)} \tag{51}
\end{aligned}
$$

where $\lambda_{\min}(X)$ and $\lambda_{\max}(X)$ denote the minimum and maximum eigenvalues of $X$. Now, $\tilde{B}\tilde{B}^\dagger$ is a complex Wishart [40] matrix $\tilde{W}(N, N)$. The distribution of the minimum eigenvalue of this matrix is given by [23]

$$
p_{\lambda_{\min}}(\lambda) = \frac{N}{2} e^{-\lambda N/2}. \tag{52}
$$

Using (52) and (51), we obtain the required result. For $N = M$, the inequality in (51) can be replaced by equality if the scalar multiplier of $\lambda_{\min}$ in the denominator is replaced by $\left(1 + \frac{M}{\rho T}\right)$.

## APPENDIX B
## PROOF OF LEMMA 3

In this section we will use the notations of Lemma 1. Let $K$ be defined such that

$$
\left(I_N + L^\dagger L\right)^{1/2} = I_N + K.
$$

Using this definition of $K$, one can express (7) as

$$
F = \tilde{Q}^\dagger E - K(\tilde{R} + \tilde{Q}^\dagger E).
$$

From (8) we have

$$
\begin{aligned}
W &= (E - \tilde{Q}F)(\tilde{R} + F)^{-1} \\
&= \left(I_T - \tilde{Q}\tilde{Q}^\dagger\right) E \left(\tilde{R} + \tilde{Q}^\dagger E\right)^{-1} (I_N + K)^{-1} \\
&\quad - \tilde{Q}\left(I_N - (I_N + K)^{-1}\right) \\
&= Q_2^\perp \left(Q_2^\perp\right)^\dagger E(\tilde{R} + \tilde{Q}^\dagger E)^{-1}(I_N + K)^{-1} \\
&\quad - \tilde{Q}\left(I_N - (I_N + K)^{-1}\right) \\
&= Q_2^\perp L \left(I_N + L^\dagger L\right)^{-1/2} - \tilde{Q} + \tilde{Q}\left(I_N + L^\dagger L\right)^{-1/2} \\
&= \left(Q_2^\perp L + \tilde{Q}\right)\left(I_N + L^\dagger L\right)^{-1/2} - \tilde{Q}.
\end{aligned}
$$

Denoting the SVD of $L$ by $U_L \Sigma_L V_L^\dagger$ and letting $\Xi = \arccos\left(\left(I_N + \Sigma_L^2\right)^{-1/2}\right)$, we obtain the desired form of $W$.

## APPENDIX C
## PROOF OF THEOREM 1

Using Lemma 3 and the definition of the chordal Frobenius norm in (19), one can verify that

$$
\|W\| = d(Q_Y, Q_X) = 2\left\|\sin\left(\frac{1}{2}\Xi\right)\right\|
$$

where $\Xi$ is defined as in Lemma 3. We now show that the normal and tangential components of the additive perturbation term $W$ in (12) are independent of $Q_X$. Using (12), for $N = M$ we have that

$$
W = [Q_X \quad Q_X^\perp]\begin{bmatrix} -2V_L \sin^2\left(\frac{1}{2}\Xi\right) V_L^\dagger \\ 2U_L \sin\left(\frac{1}{2}\Xi\right)\cos(\frac{1}{2}\Xi)V_L^\dagger \end{bmatrix}.
$$

Therefore, the tangential component is given by

$$
Q_X^\dagger W = -2V_L \sin^2\left(\frac{1}{2}\Xi\right) V_L^\dagger \tag{53}
$$

and the normal component is given by

$$
\left(Q_X^\perp\right)^\dagger W = 2U_L \sin\left(\frac{1}{2}\Xi\right)\cos\left(\frac{1}{2}\Xi\right)V_L^\dagger. \tag{54}
$$

These components are functions of $Q_X$ through the SVD of the matrix $L$, which is, in turn, a random quantity that only depends on $G$ and $\hat{G}$ in (9); cf. (14). The result in Lemma 7 (see

Appendix A) implies that $G$ and $\hat{G}$ are independent of $Q_X$, and hence we conclude that the tangential and normal components (and hence the norm) of $W$ are independent of $Q_X$ as claimed.

## APPENDIX D
## PROOF OF PROPOSITION 1

In order to prove this proposition, we begin by stating the following result.

*Proposition 2:* For $T = 2M$, if $Q_{X_1}$, $Q_{X_1}^\perp$ and $Q_{X_2}$ belong to the constellation $\mathcal{C}$ generated by (23), then augmenting the constellation by $Q_{X_2}^\perp$ does not reduce the minimum distance between constellation points. $\square$

To prove this proposition, we merely need to show that the constellation point $Q_{X_1}^\perp$ does not increase the minimum chordal Frobenius norm between constellation points if those points exist as pairs of the form $\left( Q_{X_0}, Q_{X_0}^\perp \right)$, where $Q_{X_0}$ and $Q_{X_0}^\perp$ span orthogonal subspaces. Alternatively, we need to show that for $d(Q_{X_i}, Q_{X_j})$ defined as in (20)

$$d(Q_{X_0}, Q_{X_1}) = d\left( Q_{X_0}^\perp, Q_{X_1}^\perp \right). \tag{55}$$

We will use the following lemma.

*Lemma 8:* Let $Q_A$ and $Q_B$ be $T \times M$ unitary matrices with $T \geq 2M$. Let $Q_A^\perp$ be the unitary matrix whose columns span the null space of $Q_A$. If $U \Sigma V^\dagger$ denotes the SVD of $Q_A^\dagger Q_B$, then the SVD of $(Q_A^\perp)^\dagger Q_B$ is given by $\hat{U} \left[ \left( I_M - \Sigma^2 \right)^{1/2} \quad 0 \right]^\dagger V^\dagger$ for some unitary matrix $\hat{U}$, where 0 denotes the all zero matrix of dimension $M \times (T - 2M)$. $\square$

*Proof:* Observe that $Q_A^\perp (Q_A^\perp)^\dagger$ is a projector on the null space of $Q_A$. The uniqueness of the projector operator [25] implies that $Q_A^\perp (Q_A^\perp)^\dagger = I_T - Q_A Q_A^\dagger$. Hence

$$Q_B^\dagger Q_A^\perp \left( Q_A^\perp \right)^\dagger Q_B = I_M - Q_B^\dagger Q_A Q_A^\dagger Q_B$$
$$= I_M - V \Sigma^2 V^\dagger$$
$$= V \left( I_M - \Sigma^2 \right) V^\dagger.$$

Using the Cholesky factorization, we conclude that the SVD of $Q_B^\dagger Q_A^\perp$ is given by $V [ (I_M - \Sigma^2)^{1/2} \quad 0 ] \hat{U}^\dagger$, for some unitary matrix $\hat{U}$. $\square$

By applying Lemma 8 twice; first with $Q_A = Q_{X_0}$ and $Q_B = X_1$, and then with $Q_A = Q_{X_1}$ and $Q_B = Q_{X_0}^\perp$, we obtain the desired result.

Using the result in Proposition 2, we now prove Proposition 1 by induction.

*Proof:* For $K = 1$, $|\mathcal{C}| = 2$, and the constellation generated by (23) is given by $\mathcal{C} = \{ Q_{X_1}, Q_{X_1}^\perp \}$. Suppose that the assumption is true for $K = k$. That is, $\mathcal{C} = \{ Q_{X_1}, Q_{X_1}^\perp, \ldots, Q_{X_k}, Q_{X_k}^\perp \}$. We wish to prove that the property holds for $K = k + 1$. Let $Q_{X_{k+1}}$ be the $(k+1)$th constellation point generated by (23) with $|\mathcal{C}| = 2(k+1)$. Then by Proposition 2, we have $Q_{X_{k+1}}^\perp \in \mathcal{C}$. The proof is complete. $\square$

## APPENDIX E
## PROOF OF THEOREM 2

In order to prove the first statement of the theorem (cf. (38)), we have from (7) that $R_Y$ corresponds to a rotated version of $(\tilde{R} + F)$. Therefore, we can write

$$R_Y = \Phi(I + LL^\dagger)^{1/2} B^{-1}$$

where $L = \hat{G}_2 B^{-1}$ (cf. (14), (9)), $B$ is defined in (11) and $\Phi$ is some unitary matrix that specifies the rotation of the basis of the subspace spanned by the columns of $Q_Y$. The result of Lemma 7 in Appendix A shows that $B$ and $\hat{G}_2$ are statistically independent of $Q_X$, and hence $R_Y$ is also statistically independent of $Q_X$.

In order to prove the second statement (cf. (39)), let $\Psi = \begin{bmatrix} \cos(\Xi(H)) V_L^\dagger \\ \sin(\Xi(H)) V_L^\dagger \end{bmatrix}$ and $U_X = \begin{bmatrix} Q_X & Q_X^\perp \end{bmatrix}$, where $\Xi$ is defined in Lemma 3. The matrix $Q_X$ and hence the matrix $U_X$ are isotropically distributed. Then, using Lemmas 1 and 3, one can write

$$Q_Y = U_X \Psi.$$

In order to show that $I(Q_Y; H) = 0$, one can equivalently show that $h(Q_Y | H) = h(Q_Y)$. To that end, we have

$$h(Q_Y) \geq h(Q_Y | H) \geq h(Q_Y | \Psi). \tag{56}$$

The first inequality in (56) follows from the fact that conditioning reduces entropy [24]. The second inequality in (56) follows from the fact that $H$ affects $Q_Y$ only through $\Psi$. That is, $\Psi$ contains more information about $Q_Y$ than $H$ because it involves information about the noise components.

We now prove that $p(Q_Y | \Psi) = p(Q_Y)$. In order to do that, we only need to show that

$$p(U_X \Psi | \Psi) = g(U_X) \tag{57}$$

for some function $g$. Let $U_\Psi = \begin{bmatrix} \Psi & \Psi^\perp \end{bmatrix}$ be a square unitary matrix, where $\Psi^\perp$ is the orthogonal complement of $\Psi$. Now, because $U_\Psi$ contains more information about $Q_Y$ than $\Psi$, we have

$$h(Q_Y | \Psi) \geq h(Q_Y | U_\Psi). \tag{58}$$

In order to evaluate the differential entropy $h(Q_Y | U_\Psi)$, we consider the marginal distribution $p(Q_Y | U_\Psi)$, which is given by

$$p(Q_Y | U_\Psi) = p \left( U_X U_\Psi \begin{bmatrix} I_M \\ 0_{T - M \times M} \end{bmatrix} \Big| U_\Psi \right)$$
$$= \int p(U_X U_\Psi | U_\Psi) d(U_X \Psi^\perp). \tag{59}$$

Consider the marginal probability distribution $p(U_X U_\Psi | U_\Psi)$. Since for a fixed unitary matrix $U_\Psi$, the matrix $U_X U_\Psi$ is unitary and isotropically distributed, we have [6]

$$p(U_X U_\Psi | U_\Psi) = \alpha \int d\Omega e^{i \text{Tr}(\Omega(U_\Psi^\dagger U_X^\dagger U_X U_\Psi - I_T))}$$
$$= \alpha \int d\Omega e^{i \text{Tr}(U_\Psi \Omega U_\Psi^\dagger (U_X^\dagger U_X - I_T))} \tag{60}$$

where the integration is over the space of Hermitian matrices $\Omega$ and $\alpha$ is a normalizing scalar. We now perform a change of variables with $\Upsilon$ being the Hermitian matrix $\Upsilon = U_\Psi \Omega U_\Psi^\dagger$. That is, $\Omega = U_\Psi^\dagger \Upsilon U_\Psi$. Using a result in [40], it can be shown that the differential

$$d\Upsilon = d\Omega.$$

Using this result with the expression in (60), we have

$$p(U_X U_\Psi | U_\Psi) = \alpha \int d\Upsilon e^{i \mathrm{Tr}(\Upsilon(U_X^\dagger U_X - I_T))}. \qquad (61)$$

The right hand side of (61) is not a function of $U_\Psi$. Hence, $p(Q_Y | U_\Psi) = p(Q_Y)$. Using this result in (56) and (58) completes the proof.

## APPENDIX F
### THRESHOLD VALUES FOR THE REDUCED SEARCH DETECTOR

In this section, we propose a method for determining the threshold values $A_Y$ and $B_Y$ for the reduced search quasi-ML detection strategy described in Section V-B. As pointed out in that section, proper selection of these values is critical in controlling the tradeoff between complexity and performance of the reduced search detection strategy. Our method for selecting $A_Y$ and $B_Y$ is based on using Chebychev's inequality to bound the probability of missing the correct constellation point in the set of candidate constellation points. In the initial development we will first focus on the $N = M$ case. Later, we will discuss the extension to the case when $N > M$.

*$N = M$ Case:* Let $P_m$ be the probability of missing the correct symbol $Q_X$ in the set of candidate points. That is, $P_m$ denotes the probability that $Q_X \notin \mathcal{C}'$, where $Q_X$ is the transmitted symbol, $\mathcal{C}' \triangleq \{Q_X \in \mathcal{C} | A_Y \le d_D(Q_X, Q_{\mathrm{ref},1}) - d_D(Q_Y, Q_{\mathrm{ref},1}) < B_Y\}$ and $d_D(\cdot, \cdot)$ is the metric used by the reduced search detector. Although the chordal Frobenius norm might be more appropriate for computing the threshold values, the analytical derivation appears to be intractable when th at norm is used. Therefore, for convenience we will choose the squared projection Frobenius norm for computing the threshold values; cf. (21). That is

$$d_D(Q_{X_2}, Q_{X_1}) = \frac{1}{2} \left\| Q_{X_2} Q_{X_2}^\dagger - Q_{X_1} Q_{X_1}^\dagger \right\|_F^2.$$

Define $\xi$ to be the following function of the transmitted signal matrix $Q_X$, the received signal matrix $Y = Q_X H + \sqrt{\frac{M}{\rho T}} V = Q_Y R_Y$ and the reference point $Q_{\mathrm{ref},1}$

$$\xi = \frac{1}{2} \| Q_X Q_X^\dagger - Q_{\mathrm{ref},1} Q_{\mathrm{ref},1}^\dagger \|^2$$
$$- \frac{1}{2} \| Q_Y Q_Y^\dagger - Q_{\mathrm{ref},1} Q_{\mathrm{ref},1}^\dagger \|^2$$

$$= \mathrm{Tr}\left( Q_{\mathrm{ref},1}^\dagger (W + Q_X)(W + Q_X)^\dagger Q_{\mathrm{ref},1} \right)$$
$$- \mathrm{Tr}\left( Q_{\mathrm{ref},1}^\dagger Q_X Q_X^\dagger Q_{\mathrm{ref},1} \right) \qquad (62)$$

where we used Lemma 1 in arriving at (62). Let $\bar{\xi} = \mathrm{E}\{\xi | Y\}$ and $\sigma_\xi^2 = \mathrm{E}\{(\xi - \bar{\xi})^2 | Y\}$, where the expectation is taken with respect to the isotropically distributed transmitted signal $Q_X$, the channel $H$ and noise $V$. For a given width of the search region, $\tau > 0$, the probability of missing $P_m$ can be expressed as

$$P_m = P\left( |\xi - \bar{\xi}| > \tau \right). \qquad (63)$$

Direct computation of the probability in (63) is quite complicated and depends on the choice of the reference point. Using Chebyshev's inequality, we have, for any $\tau > 0$

$$P(|\xi - \bar{\xi}| \ge \tau) \le \frac{\sigma_\xi^2}{\tau^2}. \qquad (64)$$

This inequality suggests that for $P_m$ to be bounded by a certain constant, a good choice of $\tau$ scales with $\sigma_\xi$, i.e.

$$\tau = k \sigma_\xi. \qquad (65)$$

(We will discuss later in this section how to choose the value of $k$.) The corresponding threshold values, $A_Y$ and $B_Y$, are given by

$$A_Y = -\tau + \bar{\xi}, \quad B_Y = \tau + \bar{\xi}. \qquad (66)$$

In order to determine the values of $A_Y$ and $B_Y$, we need to find $\bar{\xi}$ and $\sigma_\xi$. To do so, we will first show that $\bar{\xi}$ and $\sigma_\xi$ are independent of the actual choice of the reference point. We will then proceed to compute $\bar{\xi}$ and $\sigma_\xi$. Since the transmitted signal, $Q_X$, the channel realizations, $H$, and noise, $V$, are all independent random processes, we can compute $\bar{\xi}$ and $\sigma_\xi$ by first performing the expectation over $Q_X$ and then the expectations over $H$ and $V$.

We begin by using Lemma 1 to expose the dependency of $\xi$ on $H$, $V$ and $Q_X$. From (5) we have $(W + Q_X)R_Y = Q_X H + \sqrt{\frac{M}{\rho T}} V$. Let $U_X = [Q_X \quad Q_X^\perp]$, $G = \sqrt{\frac{M}{\rho T}} Q_X^\dagger V$ and $\hat{G} = \sqrt{\frac{M}{\rho T}} (Q_X^\perp)^\dagger V$; cf. (9). Substituting into (62) we obtain (67), which appear at the bottom of the page, where $Q'_{\mathrm{ref},1} = U_X^\dagger Q_{\mathrm{ref},1}$. Observe that the averaging over $Q_X$ has been absorbed into the averaging over the reference point. This indicates that as far as the moments of $\xi$ are concerned, the reference points are identical.

The noise matrices $G$, $\hat{G}$ and the channel matrix $H$ are independent Gaussian random matrices. In particular, the elements

$$\xi = \mathrm{Tr}\left( Q_{\mathrm{ref},1}^\dagger U_X U_X^\dagger (Q_X H + E) R_Y^{-1} R_Y^{-\dagger} (Q_X H + E)^\dagger U_X U_X^\dagger Q_{\mathrm{ref},1} \right) - \mathrm{Tr}\left( Q_{\mathrm{ref},1}^\dagger U_X U_X^\dagger Q_X Q_X^\dagger U_X U_X^\dagger Q_{\mathrm{ref},1} \right)$$
$$= \mathrm{Tr}\left( (Q'_{\mathrm{ref},1})^\dagger \begin{bmatrix} H + G \\ \hat{G} \end{bmatrix} R_Y^{-1} R_Y^{-\dagger} [H^\dagger + G^\dagger \quad \hat{G}^\dagger] Q'_{\mathrm{ref},1} \right) - \mathrm{Tr}\left( (Q'_{\mathrm{ref},1})^\dagger \begin{bmatrix} I_M & 0 \\ 0 & 0 \end{bmatrix} Q'_{\mathrm{ref},1} \right) \qquad (67)$$

of $G \sim \mathcal{CN}\left(0, \frac{M}{\rho T}\right)$, and those of $\hat{G} \sim \mathcal{CN}\left(0, \frac{M}{\rho T}\right)$, $H \sim \mathcal{CN}(0,1)$ and $H + G \sim \mathcal{CN}\left(0, \left(1 + \frac{M}{\rho T}\right)\right)$. Observe that $G$ is the component of noise that adds directly to the channel. That is, $G$ contributes additively to the received signal power. However, at high SNR, the variance of this component decays as $1/\rho$, and hence the effect becomes negligible.

Assuming that the channel matrix $H$ is nonsingular, from Lemma 1, we have that at high SNR

$$
\begin{aligned}
R_Y &= \Phi\left(I + (H+G)^{-\dagger}\hat{G}^{\dagger}\hat{G}(H+G)^{-1}\right)^{-1/2}(H+G) \\
&\approx \Phi\left(I + H^{-\dagger}\hat{G}^{\dagger}\hat{G}H^{-1}\right)^{-1/2}H \\
&= \Phi\left(H^{-\dagger}\left(H^{\dagger}H + \hat{G}^{\dagger}\hat{G}\right)H^{-1}\right)^{-1/2}H \\
&\approx \Phi H \qquad\qquad\qquad\qquad\qquad\qquad\qquad (68)
\end{aligned}
$$

where we have used the fact that $G$ and $\hat{G}$ have zero mean, and the variance of their elements decays as $1/\rho$. Hence, $H + G$ and $H^{\dagger}H + \hat{G}^{\dagger}\hat{G}$ converge in distribution to $H$ and $H^{\dagger}H$, respectively, as $\rho \to \infty$. The matrix $\Phi$ is a unitary matrix that specifies the orientation of the basis within the subspace.[5] Substituting (68) into (67), we obtain

$$
\xi \approx \mathrm{Tr}\left((Q'_{\mathrm{ref},1})^{\dagger}\begin{bmatrix} 0 & R_Y^{-\dagger}\hat{G}^{\dagger} \\ \hat{G}R_Y^{-1} & \hat{G}R_Y^{-1}R_Y^{-\dagger}\hat{G}^{\dagger} \end{bmatrix}Q'_{\mathrm{ref},1}\right) \quad (69)
$$

where the approximation in (69) is obtained by assuming that the SNR is sufficiently high so that $H + G$ in (67) is approximately equal to $H$. The approximation error becomes negligible as $\rho \to \infty$.

Using the result in (68), we can assume that $R_Y$ is independent of $Q'_{\mathrm{ref},1}$ and $\hat{G}$, which enables us to compute $\bar{\xi}$ and $\sigma_{\xi}^2$ based on the approximate expression in (69). To begin with, we will compute the conditional expectations $\bar{\xi}_{|\hat{G}}$ and $\sigma_{\xi|\hat{G}}$ and then we will compute the expectation over $\hat{G}$. If we let

$$
P = \begin{bmatrix} 0 & R_Y^{-\dagger}\hat{G}^{\dagger} \\ \hat{G}R_Y^{-1} & \hat{G}R_Y^{-1}R_Y^{-\dagger}\hat{G}^{\dagger} \end{bmatrix} \qquad (70)
$$

then

$$
\begin{aligned}
\bar{\xi}_{|\hat{G}} &= \mathrm{E}\left\{\mathrm{Tr}\left(Q'_{\mathrm{ref},1}(Q'_{\mathrm{ref},1})^{\dagger}P\right)\right\} \\
&= \frac{M}{T}\mathrm{Tr}\left(\hat{G}R_Y^{-1}R_Y^{-\dagger}\hat{G}^{\dagger}\right) \qquad\qquad (71)
\end{aligned}
$$

where we have used the fact that for isotropically distributed $Q'_{\mathrm{ref},1}$, a result in [41] implies that

$$
\mathrm{E}\left\{Q'_{\mathrm{ref},1}(Q'_{\mathrm{ref},1})^{\dagger}\right\} = \frac{M}{T}I_T.
$$

In order to compute $\sigma_{\xi|\hat{G}}$, we begin by finding

$$
\begin{aligned}
\mathrm{E}\{\xi^2|\hat{G}\} &= \mathrm{E}\left\{\mathrm{Tr}^2\left((Q'_{\mathrm{ref},1})^{\dagger}PQ'_{\mathrm{ref},1}\right)\right\} \\
&= \sum_{j,k,m,n=1}^{T} p_{jk}p_{mn}\sum_{i,\ell=1}^{M}\mathrm{E}\left\{q_{ji}^*q_{ki}q_{ml}^*q_{nm}\right\} \quad (72)
\end{aligned}
$$

where $q_{ij}$ and $p_{ij}$ are $ij$th entries of $Q'_{\mathrm{ref},1}$ and $P$, respectively. Notice that our result here complements the result in [5] in which the variance of the same quadratic term is computed for Gaussian matrices rather than the isotropically distributed unitary matrices. In order to compute the statistical expectation in (72), we make use of the following lemma from [41], which we have extended to the case of complex unitary matrices.

*Lemma 9:* Let $Q$ be a $T \times M$ random isotropically distributed unitary matrix. The following statements hold.

  i) The multiplication of a fixed row or column of $Q$ by $e^{i\theta}$ does not destroy the isotropic distribution. Consequently the mixed moments of elements of a random unitary matrix are zero unless each index occurs an even number of times in the product of which we take expectation.

  ii) For all $i \in [1,T]$ and $j \in [1,M]$, $\alpha_1 \triangleq \mathrm{E}\{|q_{ij}|^4\} = \frac{2}{T(T+1)}$.

  iii) For all $i \in [1,T]$, $j \in [1,M]$ and $k \in [1,M], k \neq j$, $\alpha_2 \triangleq \mathrm{E}\{|q_{ij}|^2|q_{ik}|^2\} = \frac{1}{T(T+1)}$.

  iv) For all $i \in [1,T]$, $\ell \in [1,T], \ell \neq i$ and $j \in [1,M]$, $\mathrm{E}\{|q_{ij}|^2|q_{\ell j}|^2\} = \frac{1}{T(T+1)} = \alpha_2$.

  v) For all $i \in [1,T]$, $\ell \in [1,T], \ell \neq i$, $j \in [1,M]$ and $k \in [1,M], k \neq j$, $\alpha_3 \triangleq \mathrm{E}\{|q_{ij}|^2|q_{\ell k}|^2\} = \frac{1}{(T-1)(T+1)}$.

  vi) For all $i \in [1,T]$, $\ell \in [1,T], \ell \neq i$, $j \in [1,M]$ and $k \in [1,M], k \neq j$, $\alpha_4 \triangleq \mathrm{E}\{q_{ij}^*q_{ik}q_{\ell j}q_{\ell k}^*\} = \frac{-1}{(T-1)T(T+1)}$. $\square$

It follows from part i) of Lemma 9, that only the terms with the following indices will be nonzero in the expectation: $\{j = k = m, i = \ell\}$, $\{j = k \neq m, i = \ell\}$, $\{j \neq k = m, i = \ell\}$, $\{j = k = m, i \neq \ell\}$, $\{j = k \neq m, i \neq \ell\}$, and $\{j \neq k = m, i \neq \ell\}$.

Using the results in Lemma 9, after regrouping and arranging terms, we can express $\mathrm{E}\{\xi^2|\hat{G}\}$ as

$$
\begin{aligned}
\mathrm{E}\{\xi^2|\hat{G}\} &= (\alpha_2 M + \alpha_4 M(M-1))\mathrm{Tr}(P^2) \\
&\quad + (\alpha_2 M + \alpha_3 M(M-1))\left(\mathrm{Tr}(P)\right)^2. \quad (73)
\end{aligned}
$$

In order to compute $\bar{\xi}$ and $\sigma_{\xi}$, we still need to average the expressions in (71) and (73) over $\hat{G}$. To that end, we observe that $\hat{G}$ is an isotropically distributed Gaussian random matrix. Hence, using the SVD of $R_Y^{-1} = U_Y \Sigma_Y V_Y^{\dagger}$

$$
\bar{\xi}_{|\hat{G}} = \frac{M}{T}\mathrm{Tr}(\hat{G}R_Y^{-1}R_Y^{-\dagger}\hat{G}^{\dagger}) = \frac{M}{T}\sum_{i=1}^{M}\sigma_i^2\|g_i\|^2
$$

where we have used $g_i$ to denote the $i$th column of $\hat{G}$ and $\sigma_i$ to denote the $i$th diagonal entry of $\Sigma_Y$. We have also used the fact that $\hat{G}U_Y \overset{d}{=} \hat{G}$, where $\overset{d}{=}$ denotes equality in distribution.

$$
\begin{aligned}
\mathrm{Tr}(P^2) &\overset{d}{=} 2\mathrm{Tr}\left(\hat{G}\Sigma_Y^2\hat{G}^{\dagger}\right) + \mathrm{Tr}\left(\Sigma_Y\hat{G}^{\dagger}\hat{G}\Sigma_Y^2\hat{G}^{\dagger}\hat{G}\Sigma_Y\right) \\
&= 2\sum_{i=1}^{M}\sigma_i^2\|g_i\|^2 + \|\Sigma_Y\hat{G}^{\dagger}\hat{G}\Sigma_Y\|^2 \\
&= 2\sum_{i=1}^{M}\sigma_i^2\|g_i\|^2 + \sum_{i=1}^{M}\sigma_i^4\|g_i\|^4 \\
&\quad + \sum_{i=1}^{M}\sum_{\substack{j=1 \\ j\neq i}}^{M}\sigma_i^2\sigma_j^2\left|g_i^{\dagger}g_j\right|^2.
\end{aligned}
$$

---

[5]Note that the receiver has no access to $\Phi$. That is, the receiver does not know the rotation of the QR decomposition that corresponds to the signal propagation through the channel.

Likewise

$$\left(\mathrm{Tr}(P)\right)^2 \stackrel{d}{=} \sum_{i=1}^{M} \sigma_i^4 \|g_i\|^4 + \sum_{i=1}^{M} \sum_{\substack{j=1 \\ j \neq i}}^{M} \sigma_i^2 \sigma_j^2 \|g_i\|^2 \|g_j\|^2. \quad (74)$$

Now, $\|g_i\|^2$ is a Chi-square random variable with $2(T - M)$ degrees of freedom and Gaussian variance equal to $\frac{M}{2\rho T}$. Hence

$$\bar{\xi} = \frac{M^2(T - M)}{T^2 \rho} \sum_{i=1}^{M} \sigma_i^2 \quad (75)$$

$$\mathrm{E}\{\mathrm{Tr}(P^2)\} = 2\frac{M(T - M)}{T\rho} \sum_{i=1}^{M} \sigma_i^2$$
$$+ \frac{M^2(T - M)(T - M + 1)}{T^2 \rho^2} \sum_{i=1}^{M} \sigma_i^4$$
$$+ \frac{M^2(T - M)}{T^2 \rho^2} \sum_{i=1}^{M} \sum_{\substack{j=1 \\ j \neq i}}^{M} \sigma_i^2 \sigma_j^2 \quad (76)$$

$$\mathrm{E}\{\left(\mathrm{Tr}(P)\right)^2\} = \frac{M^2(T - M)(T - M + 1)}{T^2 \rho^2} \sum_{i=1}^{M} \sigma_i^4$$
$$+ \frac{M^2(T - M)^2}{T^2 \rho^2} \sum_{i=1}^{M} \sum_{\substack{j=1 \\ j \neq i}}^{M} \sigma_i^2 \sigma_j^2. \quad (77)$$

Using (76), (77), (73), (75) and the fact that

$$\sigma_\xi^2 = \mathrm{E}\{\mathrm{E}\{\xi^2 | \hat{G}\}\} - \bar{\xi}^2$$

one can readily compute the proposed values of $A_Y$ and $B_Y$ given in (66) for any given $R_Y$.

We now make a few remarks regarding the values computed for $A_Y$ and $B_Y$.

- The high SNR asymptotic result in (68) conforms with the result of Theorem 2 that the channel information is captured by the $R$-component of the received signal. However, Theorem 2 asserts a stronger claim; at any SNR the channel information is contained in $R_Y$.
- Notice the dependence of $A_Y$ and $B_Y$ on the the singular values of $R_Y$. Specifically, as the singular values of $R_Y^{-1}$ grow, indicating that the channel is close to singularity (cf. (68)), the width of the search region should be increased in order to maintain a reasonable likelihood of the transmitted constellation point being within the search region.
- Notice that, in agreement with Theorem 2, the values of $A_Y$ and $B_Y$ turn out to be independent of the received signal subspace spanned by $Q_Y$.
- As the SNR, $\rho \to \infty$, $\sigma_\xi^2$ decays as $\frac{1}{\rho}$. Therefore, $P_m \to 0$ (cf. (64)), and, for a fixed $k$, the width of the search region approaches zero as $\frac{1}{\sqrt{\rho}}$. In that case, if more than one reference point is used, only a small number of likelihoods need to be computed. This observation is verified numerically in Section VII.

It remains to determine an appropriate choice for the parameter $k$ in (65). To that end we have the following comments.

i) The above analysis was based on using one reference point. When increasing the number of reference points, we found it beneficial to slightly increase the scaling parameter (65). Since the distribution of $\xi$ in (67) seems complicated to compute, even in the case of one reference point, the computation of an exact scaling parameter does not seem feasible. We will therefore assume that $\xi$ is a zero-mean Gaussian random variable, and determine a method for adjusting the scaling factor according to the number of reference points, $N_{\mathrm{ref}}$. Suppose that $k_1$ in (65) is the (positive real-valued) parameter used when one reference point is used and we wish to find the corresponding value $k_{N_{\mathrm{ref}}}$ associated with a choice of $N_{\mathrm{ref}}$ different reference points. For $N_{\mathrm{ref}}$ points, the event of missing occurs if the correct constellation point does not lie in the prescribed band associated with the first reference point or it does not lie in the prescribed band associated with the second reference point, and so on. Since our estimation of the missing probability is based on the moments of $\xi$ (cf. (64)), which are invariant under the choice of the reference point, we conclude that the estimate of the probability of missing is simply the sum of probabilities that the correct constellation point does not lie in one of the prescribed bands. That is, $P_m = N_{\mathrm{ref}} P_{m_0}$, where $P_{m_0}$ is the probability of missing with respect to the first reference point, $Q_{\mathrm{ref},1}$. Our goal is to find $k_{N_{\mathrm{ref}}}$ such that

$$N_{\mathrm{ref}} P_{m_0}(k_{N_{\mathrm{ref}}}) = P_{m_0}(k_1). \quad (78)$$

Under the Gaussian assumption on $\xi$, and with the choice of $\tau$ in (65), we have from (78), that

$$k_{N_{\mathrm{ref}}} = \sqrt{2}\mathrm{erfc}^{-1}\left(\frac{1}{N_{\mathrm{ref}}}\mathrm{erfc}(k_1/\sqrt{2})\right). \quad (79)$$

This choice of $k_{N_{\mathrm{ref}}}$, though approximate, provides reasonable guidance as to how $k$ should be adjusted with the number of reference points.

ii) It is well known that the bound given by Chebychev's inequality (64) is quite loose. A better bound can be derived using the Chernoff bound. However, the latter involves computations that require the distribution of the entries of random unitary matrices. One approach might be to assume Gaussianity of the entries, which is an asymptotically tight assumption [41], but this is beyond our current scope.

iii) We now consider the effect of choosing a particular $k$ on the probability of error and the detection complexity. Let the probability of error be denoted as $P_e(\rho)$. The probability of error is the probability that the constellation point in $\mathcal{C}'$ with the highest likelihood, say $\hat{Q}_X$, is not the transmitted constellation point, $Q_X$. This probability can be written as,

$$P_e(\rho) = P(\hat{Q}_X \neq Q_X | Q_X \in \mathcal{C}') P(Q_X \in \mathcal{C}')$$
$$+ P(Q_X \notin \mathcal{C}')$$
$$= (1 - P_m) P(\hat{Q}_X \neq Q_X | Q_X \in \mathcal{C}') + P_m$$
$$\approx P(\hat{Q}_X \neq Q_X | Q_X \in \mathcal{C}') + P_m \quad (80)$$

where the approximation in (80) is based on the assump-

tion that $P_m \ll 1$. The first term on the right-hand side of (80) denotes the inherent probability of error due to ML detection, whereas the second term is the probability of error due to the correct constellation point falling outside the reduced search region $\mathcal{C}'$. In order to ensure that our choice of $k$ does not result in significant deterioration in the performance of the receiver, $P_m$ must be maintained small with respect to $P(\hat{Q}_X \neq Q_X | Q_X \in \mathcal{C}')$. We note that for a given value of $k$, the Chebychev bound in (64) asserts that $P_m$ is upper bounded by $1/k^2$. However, this bound is generally loose [42] and one can afford to use values of $k$ smaller than those predicted by this bound without significantly affecting the overall probability of error. At high SNR the probability of missing $P_m$ can be large compared to the inherent probability of error of ML detection. This suggests that in order to ensure that the performance of the reduced search detector continues to be essentially the same as that of the ML detector, larger values of $k$ must be chosen at higher SNRs. That is, we should allow the scaling parameter $k$ to be a function of the SNR. However, we would still like to have $|\mathcal{C}'| \rightarrow 1$ as $\rho \rightarrow \infty$. Since $\sigma_\xi$ decays as fast as $\frac{1}{\sqrt{\rho}}$, we seek a function that grows no faster than $\sqrt{\rho}$. One such function is $c \log(\cdot)$, for some constant $c$, which results in an overall decay in the band width at a rate of $\frac{\log(\rho)}{\sqrt{\rho}}$. In Section VII we show that this choice of $k$ yields good tradeoff in performance versus decoding complexity. Hence, for any SNR and number of reference points, $N_{\mathrm{ref}}$, we propose to choose the parameter $k$ to be

$$k = c\sqrt{2}\mathrm{erfc}^{-1}\Big(\frac{1}{N_{\mathrm{ref}}}\mathrm{erfc}(k_1/\sqrt{2})\Big)\log(\rho) \qquad (81)$$

for some scalars $k_1$ and $c$, where we have used the expression in (79).

$N > M$ *Case:* In this case, the channel matrix, $H \in \mathbb{C}^{M \times N}$ is fat. Let $H$ be partitioned as $H_1 \in \mathbb{C}^{M \times M}$ and $H_2 \in \mathbb{C}^{M \times (N-M)}$, $V \in \mathbb{C}^{T \times N}$ be partitioned as $V_1 \in \mathbb{C}^{T \times M}$ and $V_2 \in \mathbb{C}^{T \times (N-M)}$. Assume that the receiver performs the unique QR-decomposition of $Y$ such that $R_Y$ is an upper triangular matrix with real nonnegative diagonal elements. That is

$$Y = Q_Y R_Y = [\,Q_{Y_1} \quad Q_{Y_2}\,] \begin{bmatrix} R_{Y_{11}} & R_{Y_{12}} \\ 0 & R_{Y_{22}} \end{bmatrix} \qquad (82)$$

$$= Q_X [\,H_1 \quad H_2\,] + \sqrt{\frac{M}{\rho T}} [\,V_1 \quad V_2\,]. \quad (83)$$

Therefore, using (82) and (83), we have

$$Q_{Y_1} R_{Y_{11}} = Q_X H_1 + \sqrt{\frac{M}{\rho T}} V_1 \qquad (84)$$

$$Q_{Y_1} R_{Y_{12}} + Q_{Y_2} R_{Y_{22}} = Q_X H_2 + \sqrt{\frac{M}{\rho T}} V_2. \qquad (85)$$

From (84), it is clear that we could determine the search region by simply applying the previous algorithm to the first $M$ columns of $Y$. In that case, the probability of missing the correct constellation point in $\mathcal{C}'$ would be the same as that when $N =$

$M$, but the information from the additional antennas is naturally incorporated into the likelihoods; cf. (41). The performance of that approach can be significantly improved by first permuting $Y$ so that the diagonal entries of $(R_Y R_Y^\dagger)^{-1}$ are arranged in a nondecreasing order [25]. The results in Appendix F-A suggest that the width of the search region is dominated by the inverse of the the minimum eigenvalue of $R_Y R_Y^\dagger$; i.e., the maximum eigenvalue of $\left(R_Y R_Y^\dagger\right)^{-1}$. Since the diagonal entries of a positive semidefinite matrix strongly majorize its eigenvalues [43], after permuting $Y$, the width of the search region that corresponds to the upper left $M \times M$ block of $R_Y$ is typically less than the width of the search region that corresponds to blocks of smaller diagonal entries. That is, the diversity gain offered by the increase in the number of receive antennas can be used to reduce the adverse effect of atypically weak channel realizations on the detection complexity. In that way, the additional receive antennas play a role in both the selection of the candidate constellation points and in the values of the likelihoods; something that improves the tradeoff between the probability of missing and the number of likelihood computations. This tradeoff could be more tightly controlled by exploiting the information in (85). However, such an algorithm would need to compute the distance not only between constellation points and reference points but also the distance between the $(N - M)$-dimensional null spaces of the $T \times M$ constellation points and the reference points. That is, the look-up table would have to be augmented in order to include the distance to reference points from all possible $(N - M)$-dimensional null spaces of the $T \times M$ constellation points. Since there are $\binom{T-M}{N-M}$ possible null spaces associated with each $T \times M$ constellation point and $T \geq \min\{M, N\} + N = M + N$, attempting to reduce the probability of missing the correct constellation point in $\mathcal{C}'$ by exploiting the null space information contained in (85) would require an increase in memory size that would be hard to justify. Fortunately, the probability of missing can be easily controlled via the scaling factor $k$ in (66) and the number of reference points at a small storage cost.

## APPENDIX G
## PAIRWISE ERROR PROBABILITY

The pairwise error probability (PEP) $P(i \rightarrow j)$ is defined as the probability that the receiver mistakes the $i$th constellation point for the $j$th one given that the $i$th constellation point has been transmitted. Assuming that the receiver employs an ML detector[6] (cf. (35)), then the PEP is given by

$$P(i \rightarrow j) = P\Big( \mathrm{Tr}\,\Big(Y^\dagger Q_{X_i} Q_{X_i}^\dagger Y\Big) \leq \mathrm{Tr}\,\Big(Y^\dagger Q_{X_j} Q_{X_j}^\dagger Y\Big)\Big). \qquad (86)$$

Since we assume $Q_{X_i}$ to be the transmitted constellation point, then the received signal can be expressed as $Y = Q_{X_i} H + \sqrt{\frac{M}{\rho T}} V$, where $H$ is the $M \times N$ isotropically distributed (i.d.) complex Gaussian channel matrix and $V$ is the $T \times N$ i.d. complex Gaussian additive noise matrix.

[6]Notice that the performance of the reduced search detector in Section V-B is very similar to that of the ML detector.

Now consider the first term in the probability argument in (86). Let $G$ and $\hat{G}$ be defined as in (9). Then

$$
\begin{aligned}
&\mathrm{Tr}\left(Y^\dagger Q_{X_i} Q_{X_i}^\dagger Y\right) \\
&= \mathrm{Tr}\left(\left(H^\dagger + \sqrt{\tfrac{M}{\rho T}} G^\dagger\right)\left(H + \sqrt{\tfrac{M}{\rho T}} G\right)\right) \\
&= \|S\|_F^2
\end{aligned}
\tag{87}
$$

where $S = \left(H + \sqrt{\tfrac{M}{\rho T}} G\right)$ is a zero mean isotropically distributed complex Gaussian random matrix with variance $\tfrac{1}{2}\left(1 + \tfrac{M}{\rho T}\right)$ per real dimension. Similarly

$$
\begin{aligned}
&\mathrm{Tr}\left(Y^\dagger Q_{X_j} Q_{X_j}^\dagger Y\right) \\
&= \left\| H^\dagger Q_{X_i}^\dagger Q_{X_j} + \sqrt{\tfrac{M}{\rho T}} V^\dagger Q_{X_j} \right\|_F^2 \\
&= \left\| H^\dagger Q_{X_i}^\dagger Q_{X_j} + \sqrt{\tfrac{M}{\rho T}} V^\dagger \left[ Q_{X_i} Q_{X_i}^\perp \right]\left[\begin{array}{c} Q_{X_i}^\dagger \\ (Q_{X_i}^\perp)^\dagger \end{array}\right] Q_{X_j} \right\|_F^2 \\
&= \left\| S^\dagger Q_{X_i}^\dagger Q_{X_j} + \hat{G}^\dagger \left(Q_{X_i}^\perp\right)^\dagger Q_{X_j} \right\|_F^2.
\end{aligned}
\tag{88}
$$

The $(T-M) \times N$ random matrix $\hat{G}$ is zero mean isotropically distributed (i.d.) complex Gaussian with variance of $\tfrac{M}{2\rho T}$ per real dimension. From Lemma 7 in Appendix A, we have that $\hat{G}$ and $G$ are statistically independent. Using the result of Lemma 8 in Appendix D and (88), the PEP in (86) can be written by (89), which appears at the bottom of the page, where $V_{ij}\Sigma_{ij}U_{ij}^\dagger$ and $V_{ij}\left[ \left(I - \Sigma_{ij}^2\right)^{1/2} \quad 0_{M \times (T-2M)} \right] \hat{U}_{ij}^\dagger$ denote the SVD of $Q_{X_j}^\dagger Q_{X_i}$ and $Q_{X_j}^\dagger Q_{X_i}^\perp$ respectively, $S_A = U_{ij}^\dagger S$ and $S_B = \left[ I_M \quad 0_{M \times (T-2M)} \right] \hat{U}_{ij}^\dagger \hat{G}$. Notice that because $T \geq 2M$ (cf. (1)), $S_A$ and $S_B$ are independent $M \times N$ i.d. Gaussian random matrices with respective variances $\tfrac{1}{2}\left(1 + \tfrac{M}{\rho T}\right)$ and $\tfrac{M}{2\rho T}$ per real dimension. The equality in (89) follows from the invariance of the norm under unitary transformation. Now, (89) can be written as (90), which also appears at the bottom of the page. Let $\tilde{S}_A$ and $\tilde{S}_B$ be the zero mean and unit variance i.d Gaussian random matrices $\sqrt{\tfrac{2\rho T}{\rho T + M}} S_A$ and $\sqrt{\tfrac{2\rho T}{M}} S_B$, respectively. If we let $S_C$

be the i.d. random Gaussian matrix with i.i.d. zero mean unit variance entries given by $S_C = \left[ \tilde{S}_A^\dagger \quad \tilde{S}_B^\dagger \right]$, and let $\Gamma_{ij}$ be given by

$$
\Gamma_{ij} = \left[\begin{array}{cc} \gamma_{11}\left(\Sigma_{ij}^2 - I\right) & \gamma_{12}\left(I - \Sigma_{ij}^2\right)^{1/2}\Sigma_{ij} \\ \gamma_{12}\left(I - \Sigma_{ij}^2\right)^{1/2}\Sigma_{ij} & \gamma_{22}\left(I - \Sigma_{ij}^2\right) \end{array}\right]
$$

where $\gamma_{11} = \tfrac{1}{2}\left(1 + \tfrac{M}{\rho T}\right)$, $\gamma_{22} = \tfrac{1}{2}\tfrac{M}{\rho T}$, and $\gamma_{12} = \sqrt{\gamma_{11}\gamma_{22}}$, then one can express the argument of the trace in (90) as

$$
\begin{aligned}
&S_A^\dagger\left(\Sigma_{ij}^2 - I\right) S_A + S_A^\dagger \Sigma_{ij}\left(I - \Sigma_{ij}^2\right)^{1/2} S_B + S_B^\dagger \Sigma_{ij} \\
&\times \left(I - \Sigma_{ij}^2\right)^{1/2} S_A + S_B^\dagger\left(I - \Sigma_{ij}^2\right) S_B = S_C \Gamma_{ij} S_C^\dagger.
\end{aligned}
\tag{91}
$$

Notice that $\Gamma_{ij}$ is symmetric and hence using the eigen decomposition of $\Gamma_{ij}$ we obtain, $\Gamma_{ij} = Q_{ij}\Lambda_{ij}Q_{ij}^\dagger$. Since $S_C$ is i.d., then $S_C \overset{d}{=} \tilde{S}_C = S_C Q_{ij}$ and the expression in (90) reduces to

$$
P(i \to j) = P\left(0 \leq \mathrm{Tr}\left(\tilde{S}_C \Lambda_{ij} \tilde{S}_C^\dagger\right)\right)
\tag{92}
$$

$$
= P\left(0 \leq \sum_{k=1}^{2M} \lambda_k^{ij} \|S_{C_k}\|^2\right)
\tag{93}
$$

where $S_{C_k}$ is the $k$th column of $\tilde{S}_C$. Since $\|S_{C_k}\|^2$ is a Chi-square distributed random variable with $2N$ degrees of freedom, then if we denote the multiplicity of $\lambda_k^{ij}$ by $\kappa(k)$, the expression in (45) follows.

It remains to determine the eigenvalues of $\Gamma_{ij}$ in terms of $\Sigma_{ij}$. For brevity, we will suppress the subscript $ij$ in $\Gamma_{ij}$ and $\Sigma_{ij}$, and we will define $\Theta = (I - \Sigma^2)^{1/2}$. Using these notations, we obtain

$$
\Gamma = \left[\begin{array}{cc} -\gamma_{11}\Theta^2 & \gamma_{12}\Theta\Sigma \\ \gamma_{12}\Theta\Sigma & \gamma_{22}\Theta^2 \end{array}\right]
\tag{94}
$$

where $\Theta$ and $\Sigma$ are both diagonal matrices with non negative entries given by $\theta_\ell$ and $\sigma_\ell$, respectively. In order to compute $\lambda_k$ for $k \in \{1, 2, \ldots, 2M\}$, we need to find the roots of $\det(\Gamma - \lambda I)$. Assuming that $-\gamma_{11}\Theta^2 - \lambda I$ is nonsingular, an expression for those roots can be obtained by using the expression for the

$$
\begin{aligned}
P(i \to j) &= P\left(\|S\|^2 \leq \left\| V_{ij}\Sigma_{ij}U_{ij}^\dagger S + V_{ij}\left[ \left(I - \Sigma_{ij}^2\right)^{1/2} \quad 0_{M \times (T-2M)} \right] \hat{U}_{ij}^\dagger \hat{G} \right\|^2\right) \\
&= P\left(\|U_{ij}^\dagger S\|^2 \leq \left\| (\Sigma_{ij}U_{ij}^\dagger S + \left[ \left(I - \Sigma_{ij}^2\right)^{1/2} \quad 0_{M \times (T-2M)} \right] \hat{U}_{ij}^\dagger \hat{G}) \right\|^2\right) \\
&= P\left(\|S_A\|^2 \leq \left\| \Sigma_{ij}S_A + \left(I - \Sigma_{ij}^2\right)^{1/2} S_B \right\|^2\right)
\end{aligned}
\tag{89}
$$

$$
\begin{aligned}
P(i \to j) &= P\left(\mathrm{Tr}(S_A^\dagger S_A) \leq \mathrm{Tr}\left(S_A^\dagger \Sigma_{ij}^2 S_A\right) + S_B^\dagger \Sigma_{ij}\left(I - \Sigma_{ij}^2\right)^{1/2} S_A + S_B^\dagger\left(I - \Sigma_{ij}^2\right) S_B S_A^\dagger \Sigma_{ij}\left(I - \Sigma_{ij}^2\right)^{1/2} S_B\right) \\
&= P\left(0 \leq \mathrm{Tr}(S_A^\dagger\left(\Sigma_{ij}^2 - I\right) S_A + S_B^\dagger\left(I - \Sigma_{ij}^2\right) S_B) + S_A^\dagger \Sigma_{ij}\left(I - \Sigma_{ij}^2\right)^{1/2} S_B + S_B^\dagger \Sigma_{ij}\left(I - \Sigma_{ij}^2\right)^{1/2} S_A\right).
\end{aligned}
\tag{90}
$$

determinant of a block partitioned matrix [43] and the diagonal structure of $\Theta$ and $\Sigma$. Specifically, one can show that

$$
\begin{aligned}
\det(\Gamma - \lambda I) &= \det\left(-\gamma_{12}^2\Theta^4 + \frac{\lambda}{2}\Theta^2 + \lambda^2 I - \gamma_{12}^2\Theta^2\Sigma^2\right) \\
&= (-1)^M \det\left(\gamma_{12}^2\Theta^2(\Theta^2 + \Sigma^2) - \frac{\lambda}{2}\Theta^2 - \lambda^2 I\right) \\
&= (-1)^M \det\left(\gamma_{12}^2\Theta^2 - \frac{\lambda}{2}\Theta^2 - \lambda^2 I\right) \\
&= \prod_{i=1}^{M}\left(\lambda^2 + \frac{\lambda}{2}\theta_i^2 - \gamma_{12}^2\theta_i^2\right).
\end{aligned}
$$

Therefore, the eigenvalues of $\Gamma$ are given by

$$
\lambda_i = \frac{\theta_i}{4}\left(-\theta_i \pm \sqrt{\theta_i^2 + 16\gamma_{12}^2}\right)
$$

and hence one can see that $\Gamma$ has $M$ nonpositive eigenvalues.

## REFERENCES

[1] I. E. Telatar, "Capacity of multiantenna Gaussian channels," *Eur. Trans. Telecom.*, vol. 10, pp. 585–595, Nov. 1999.

[2] L. Zheng and D. N. C. Tse, "Communication on the Grassmann manifold: A geometric approach to the noncoherent multiple-antenna channel," *IEEE Trans. Inf. Theory*, vol. 48, pp. 359–383, Feb. 2002.

[3] T. L. Marzetta and B. M. Hochwald, "Capacity of a mobile multiple-antenna communication link in Rayleigh flat fading," *IEEE Trans. Inf. Theory*, vol. 45, pp. 139–157, Jan. 1999.

[4] B. M. Hochwald and T. L. Marzetta, "Unitary space-time modulation multiple-antenna communications in Rayleigh flat fading," *IEEE Trans. Inf. Theory*, vol. 46, pp. 543–564, Mar. 2000.

[5] C. Rao and B. Hassibi, "Analysis of multiple antenna wireless links at low SNR," *IEEE Trans. Inf. Theory*, vol. 50, pp. 2123–2130, Sep. 2004.

[6] B. Hassibi and T. L. Marzetta, "Multiple-antennas and isotropically random unitary inputs: The received signal density in closed form," *IEEE Trans. Inf. Theory*, vol. 48, pp. 1473–1484, Jun. 2002.

[7] B. Hassibi and B. M. Hochwald, "How much training is needed in multiple-antenna wireless links?," *IEEE Trans. Inf. Theory*, vol. 49, pp. 951–963, Apr. 2003.

[8] A. Lapidoth and S. M. Moser, "Capacity bounds via duality with applications to multiple-antenna systems on flat-fading channels," *IEEE Trans. Inf. Theory*, vol. 49, pp. 2426–2467, Oct. 2003.

[9] D. Agrawal, T. Richardson, and R. Urbanke, "Multiple-antenna signal constellations for fading channels," *IEEE Trans. Inf. Theory*, vol. 47, pp. 2618–2626, Sep. 2001.

[10] P. Dayal, M. Brehler, and M. K. Varanasi, "Leveraging coherent space-time codes for noncoherent communication via training," *IEEE Trans. Inf. Theory*, vol. 50, pp. 2058–2080, Sep. 2004.

[11] H. El-Gamal and M. O. Damen, "Noncoherent space-time coding: An algebraic perspective," *IEEE Trans. Inf. Theory*, vol. 51, pp. 2380–2390, Jul. 2005.

[12] A. Barg and D. Y. Nogin, "Bounds on packings of spheres in the Grassmann manifold," *IEEE Trans. Inf. Theory*, vol. 48, pp. 2450–2454, Sep. 2002.

[13] O. Henkel, "Sphere-packing bounds in the Grassmann and Stiefel manifolds," *IEEE Trans. Inf. Theory*, vol. 51, pp. 3445–3456, Oct. 2005.

[14] M. J. Borran, A. Sabharwal, and B. Aazhang, "On design criteria and construction of noncoherent space-time constellations," *IEEE Trans. Inf. Theory*, vol. 49, pp. 2332–2351, Oct. 2003.

[15] J.-C. Belfiore and A. M. Cipriano, "Space-time coding for noncoherent channels," in *Space-Time Wireless Systems: From Array Processing to MIMO Communications*, H. Bölcskei, D. Gesbert, C. B. Papadias, and A.-J. van der Veen, Eds. New York: Cambridge University Press, 2006, ch. 10, pp. 198–217.

[16] V. Tarokh and I.-M. Kim, "Existence and construction of noncoherent unitary space-time codes," *IEEE Trans. Inf. Theory*, vol. 48, pp. 3112–3117, Dec. 2002.

[17] B. M. Hochwald, T. L. Marzetta, T. J. Richardson, W. Sweldens, and R. Urbanke, "Systematic design of unitary space-time constellations," *IEEE Trans. Inf. Theory*, vol. 46, pp. 1962–1973, Sep. 2000.

[18] I. Kammoun and J.-C. Belfiore, "A new family of Grassmann space-time codes for non-coherent MIMO systems," *IEEE Commun. Lett.*, vol. 7, pp. 528–530, Nov. 2003.

[19] I. Kammoun, A. M. Cipriano, and J.-C. Belfiore, "Non-coherent codes over the Grassmannian," *IEEE Trans. Wireless Commun.*, vol. 6, pp. 3657–3667, Oct. 2007.

[20] M. L. McCloud, M. Brehler, and M. K. Varanasi, "Signal design and convolutional coding for noncoherent space-time communication on the block-Rayleigh-fading channel," *IEEE Trans. Inf. Theory*, vol. 48, pp. 1186–1194, May 2002.

[21] A. Edelman, T. Arias, and S. Smith, "The geometry of algorithms with orthogonality constraints," *SIAM J. Matrix Anal. Appl.*, vol. 20, no. 2, pp. 303–353, 1998.

[22] G. W. Stewart, "Perturbation bounds for the QR factorization of a matrix," *SIAM J. Num. Anal.*, vol. 14, pp. 509–518, Jun. 1977.

[23] A. Edelman, "Eigenvalues and condition numbers of random matrices," *SIAM J. Matrix Anal. Appl.*, vol. 9, pp. 543–560, Oct. 1988.

[24] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.

[25] G. H. Golub and C. F. van Loan, *Matrix Computations*, third ed. Baltimore, MD: The Johns Hopkins University Press, 1996.

[26] J. H. Conway, R. H. Hardin, and N. J. A. Sloane, "Packing lines, planes, etc.: Packings in Grassmannian spaces," *Exper. Math.*, vol. 5, no. 2, pp. 139–159, 1996.

[27] B. M. Hochwald and S. ten Brink, "Achieving near-capacity on a multiple-antenna channel," *IEEE Trans. Commun.*, vol. 51, pp. 389–399, Mar. 2003.

[28] A. J. Goldsmith and S.-G. Chua, "Adaptive coded modulation for fading channels," *IEEE Trans. Commun.*, vol. 46, pp. 595–602, May 1998.

[29] Y. Li and X. G. Xia, "Constellation mapping for space-time matrix modulation with iterative demodulation/decoding," *IEEE Trans. Commun.*, vol. 53, pp. 764–768, May 2005.

[30] I. Bahceci and T. Duman, "Trellis-coded unitary space-time modulation," *IEEE Trans. Wireless Commun.*, vol. 3, pp. 2005–2012, Nov. 2004.

[31] V. Tarokh, N. Seshadri, and A. R. Calderbank, "Space-time codes for high data rate wireless communication: Performance criterion and code construction," *IEEE Trans. Inf. Theory*, vol. 44, pp. 744–765, 1998.

[32] M. Brehler and M. K. Varanasi, "Asymptotic error probability analysis of quadratic receivers in Rayleigh-fading channels with applications to a unified analysis of coherent and noncoherent space-time receivers," *IEEE Trans. Inf. Theory*, vol. 47, pp. 2383–2399, Sep. 2001.

[33] R. Li and P. Y. Kam, "New tight bounds on the pairwise error probability for unitary space-time modulations," *IEEE Commun. Lett.*, vol. 4, pp. 289–291, Apr. 2005.

[34] J. Robinson, "The distribution of a general quadratic form in normal variates," *Austral. J. Statist.*, vol. 7, no. 3, pp. 110–114, 1965.

[35] N. Johnson and S. Kotz, *Distributions in Statistics: Continuous Univariate Distributions—2*. New York: Wiley, 1970.

[36] S. M. Alamouti, "A simple transmitter diversity scheme for wireless communications," *IEEE J. Sel. Areas Commun.*, vol. 16, pp. 1451–1458, Oct. 1998.

[37] B. Hassibi and B. M. Hochwald, "High-rate codes that are linear in space and time," *IEEE Trans. Inf. Theory*, vol. 48, pp. 1804–1824, Jul. 2002.

[38] J.-C. Belfiore, G. Rekaya, and E. Viterbo, "The Golden code: A $2 \times 2$ full-rate space-time code with nonvanishing determinants," *IEEE Trans. Inf. Theory*, vol. 4, pp. 1432–1436, Apr. 2005.

[39] G. Taricco and E. Biglieri, "Space-time decoding with imperfect channel estimation," *IEEE Trans. Wireless Commun.*, vol. 4, pp. 1874–1888, Jul. 2005.

[40] R. J. Muirhead, *Aspects of Multivariate Statistical Theory*. New York: Wiley, 1982.

[41] C. Stein, The Accuracy of the Normal Approximation to the Distribution of the Traces of Powers of Random Orthogonal Matrices Stanford University, Dept. of Statistics, Stanford, CA, 1995, Tech. Rep. 470.

[42] A. Papoulis, *Probability, Random Variables and Stochastic Processes*, 2nd ed. New York: McGraw-Hill, 1984.

[43] R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge, U.K.: Cambridge University Press, 1999.

**Ramy H. Gohary** (M'00) received the B.Eng. (Hons.) degree in electronics and communications engineering from Assiut University, Egypt, in 1996, the M.Sc. degree in communications engineering from Cairo University, Egypt, in 2000, and the Ph.D. degree in electrical engineering from McMaster University, ON, Canada, in 2006.

He received the Natural Sciences and Engineering Research Council Visiting Fellowship award in 2007 and he is currently a Visiting Fellow with the Terrestrial Wireless Systems Branch, Communications Research Centre, Industry

Canada. His research interests include analysis and design of MIMO wireless communication systems, applications of optimization and geometry in signal processing and communications, information theoretic aspects of multiuser communication systems, and applications of iterative detection, and decoding techniques in multiple antenna and multiuser systems.

**Timothy N. Davidson** (M'96) received the B.Eng. (Hons. I) degree in electronic engineering from the University of Western Australia (UWA), Perth, in 1991 and the D.Phil. degree in engineering science from the University of Oxford, U.K., in 1995.

He is currently an Associate Professor in the Department of Electrical and Computer Engineering at McMaster University, Hamilton, Ontario, Canada, where he holds the (Tier II) Canada Research Chair in Communication Systems, and is currently serving as Acting Director of McMaster's School of Computational Engineering and Science. He is also a Registered Professional Engineer in the Province of Ontario. His research interests lie in the general areas of communications, signal processing, and control. He has held research positions at the Communications Research Laboratory at McMaster University, the Adaptive Signal Processing Laboratory at UWA, and the Australian Telecommunications Research Institute at Curtin University of Technology, Perth, Western Australia.

Dr. Davidson was awarded the 1991 J. A. Wood Memorial Prize (for "the most outstanding [UWA] graduand" in the pure and applied sciences) and the 1991 Rhodes Scholarship for Western Australia. He is currently serving as an Associate Editor of the IEEE TRANSACTIONS ON SIGNAL PROCESSING, the IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS, and *Optimization and Engineering*. He has also served as an Associate Editor of the IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS II, and as a Guest Co-editor of issues of the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS and the IEEE JOURNAL ON SELECTED TOPICS IN SIGNAL PROCESSING.